

How Should We Build A Benchmark?

Revisiting 274 Code-Related Benchmarks For LLMs

Anonymous ACL submission

Abstract

Various benchmarks have been proposed to assess the performance of large language models (LLMs) in different coding scenarios. We refer to them as code-related benchmarks. However, there are no systematic guidelines by which such a benchmark should be developed to assure its quality, reliability, and reproducibility. We propose **How2Bench** comprising a 55-criteria checklist as a set of guidelines to comprehensively govern the development of code-related benchmarks. Using HOW2BENCH, we profiled 274 code-related benchmarks released within the past decade and found concerning issues. Nearly 70% of the benchmarks did not take measures for data quality assurance; over 10% did not even open source or only partially open source. Many highly cited benchmarks have loopholes, including duplicated samples, incorrect reference codes/tests/prompts, and unremoved sensitive/confidential information. Finally, we conducted a human study involving 49 participants and *revealed significant gaps in awareness* of the importance of data quality, reproducibility, and transparency. For ease of use, we provide a *printable version* of HOW2BENCH in Appendix F.

1 Introduction

“Awareness is the beginning of action;
action is the fulfillment of awareness.”
— Yangming Wang (1472 - 1529)

Recent large language models (LLMs) have shown remarkable capabilities across various domains such as software development (Chen et al., 2021a), question answering (Rogers et al., 2023), and math reasoning (Imani et al., 2023). Various *benchmarks* (Chen et al., 2021a; Jimenez et al., 2024; Austin et al., 2021; Yue et al., 2024; Du et al., 2023a) are proposed to evaluate LLMs’ effectiveness and limitations from multiple perspectives in different application scenarios.

However, *doubts regarding the quality, reliability, and transparency* of various code-related benchmarks arise. For example, a recent study pointed out that “current programming benchmarks are inadequate for assessing the actual correctness of LLM-generated code” (Liu et al., 2023a). Other accusations, including *irreproducible* (Reuel et al., 2024), *closed* data sources (Cao et al., 2024b), *low quality* (Qiu et al., 2024a; Yadav et al., 2024a), and *inadequate validation measures* (Liu et al., 2023a), were also raised, undermining the credibility of these benchmarks and thereby their subsequent evaluation results. This motivates the need for rigorous and thorough *guidelines* to govern code-related benchmark development.

In this paper, we introduce **HOW2BENCH**, a **comprehensive guideline** consisting of a 55-criteria checklist specially designed for code-related benchmarks. This checklist *covers the entire lifecycle* of benchmark development, from *design* and *construction* to *evaluation*, *analysis*, and *release* as shown in Figure 1. It underwent multiple iterations — we initiated a draft inspired by open-source software guidelines (Fogel, 2005) and classical measurement theory (Suppes et al., 1962). We refined it through iterative discussions with practitioners, leading to the finalization of these criteria. HOW2BENCH emphasizes *reliability*, *validity*, *open access*, and *reproducibility* in the benchmark development, ensuring high standards and fostering a more reliable and transparent environment.

Following HOW2BENCH, we conducted an in-depth profiling of 270+ code-related benchmarks developed over the past decade (2014 - 2024). The extent of criteria violations by the profiled benchmarks is concerning:

- Almost 70% of the benchmarks did not take any measures for *data quality assurance*;
- Over 90% did not consider *code coverage* when use passing test cases as an oracle;
- Over half of the benchmarks did not provide

the essential information (e.g., experiment setup, prompts) for *reproducibility*;

- Over 10% are *not open source* or only partially open source.

We observed that even *highly cited benchmarks have loopholes*, including duplicated samples, incorrect reference/tests, unclear displays, and unre-moved sensitive/confidential information. We also observed these loopholes can *propagate*. Over 18% of the benchmarks serve as data sources for subsequent benchmarks (Figure 8). Therefore, the data quality of benchmarks affects their credibility and likely impacts future benchmarks.

To understand the usefulness of the criteria in How2Bench, we conducted a human study involving 49 participants through questionnaires. All participants *concurred on the necessity* of having a checklist for benchmark construction to enhance quality. Nearly all participants with experience in benchmark development *acknowledged the importance of all these 55 criteria*. The study also exposed *gaps in quality awareness*: 16% of participants were unaware of the necessity for data denoising, and over 40% were not aware that the experimental setup and environment could impact the reproducibility and transparency.

Given the identified severe situation, this paper calls for attention to the quality of LLM benchmarks due to insufficient rigor and lack of systematic, reproducible evaluation workflows. We noticed that this issue not only stems from the significant effort required to construct a rigorous benchmark, but also from a lack of awareness regarding how crucial a rigorous development process is to the quality of benchmarks and the reliability of evaluations. Without being aware of and addressing the quality of code-related benchmarks (as well as other kinds of benchmarks for LLMs), the reliability and reproducibility of benchmark results remain compromised, misleading research directions and hindering meaningful progress.

Therefore, we advocate for establishing comprehensive, transparent, and enforceable guidelines for benchmark development and quality assurance, and took an initial step, introducing HOW2BENCH. This intention also echoed several prior works (McIntosh et al., 2025; Reuel et al., 2024; Holloway et al., 1999; Zairi, 2010; Francis and Holloway, 2007; Wu et al., 2025). We hope HOW2BENCH could be helpful to elevate the standards and trustworthiness of future benchmarks.

This paper makes contributions in five aspects:

- **Novelty.** We introduce HOW2BENCH, a comprehensive set of guidelines packaged as a 55-criteria checklist that covers the lifecycle of code-related benchmark development.
- **Significance.** HOW2BENCH presents the first comprehensive set of actionable guidelines for developing high-quality benchmarks, striving to create a more reliable and transparent environment. The human study also highlighted the demand for such a detailed guideline.
- **Usefulness.** HOW2BENCH serves as a guideline for practitioners before/during developing code-related benchmarks, and a checklist for evaluating existing benchmarks after their release. For ease of use, we also provide a *printable version* of HOW2BENCH on Appendix F.
- **Generalizability.** Most criteria listed in HOW2BENCH can be adopted or adapted to other benchmarks such as Question-answering, mathematical reasoning, and multi-modal benchmarks.
- **Long-term Impact.** Our statistics alert the community to the severity and prevalence of non-standard practices in benchmark development. It ultimately improves the overall quality of benchmarks due to the propagation effect among them.

2 Background

2.1 Code-related Benchmarks

Benchmarks for coding tasks like code generation (Chen et al., 2021a; Austin et al., 2021), defect detection (Just et al., 2014; Gao et al., 2023b; Liu et al., 2024c), and program repair (Jimenez et al., 2024; Risse and Böhme, 2024) are increasingly common, reflecting the growing needs for using LLMs for coding tasks. Recent studies have highlighted various issues with these benchmarks, ranging from design inconsistencies to scope and applicability limitations. For example, (Liu et al., 2023a) found that even some widely used benchmarks, such as HumanEval (Chen et al., 2021a) and MBPP (Austin et al., 2021), contains a non-trivial proportion of bugs in implementation, documentation, and test cases. Our work, in comparison, introduces a detailed guideline that *guides the benchmark development* during the entire lifecycle.

2.2 Related Studies and Surveys

Several recent surveys and empirical studies have profiled the status quo of LLM development. These

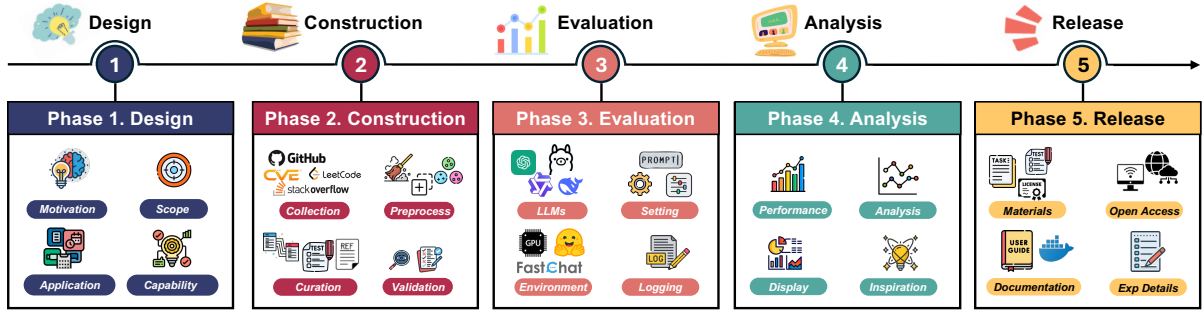


Figure 1: Lifecycle of Benchmark Development

studies either explore the overall performance for certain areas such as software engineering (Hou et al., 2023; Wang et al., 2024a) or investigate the capabilities of LLMs on specific tasks such as code generation (Dou et al., 2024; Yu et al., 2024) and test generation (Schäfer et al., 2024; Yuan et al., 2024b, 2023a). A survey (Chang et al., 2024) about how to evaluate LLMs was proposed to answer what/where/how to evaluate LLMs. This paper differs from these studies in its purpose and perspectives. Unlike these benchmarks, our work proposes guidelines for future benchmark development and provides a checklist to assess the quality of these existing benchmarks.

Recently, BetterBench (Reuel et al., 2024) is a concurrent work assessing the AI benchmarks against 46 criteria. Then, it scored 24 AI benchmarks in various domains and ranked them. BetterBench differs from this paper in several key aspects: scope (general benchmarks vs. code-related benchmarks), lifecycle division (it addresses benchmark retirement, while How2Bench focuses on benchmark evaluation, analysis, and release), and objectives (scoring benchmarks vs. offering comprehensive guidelines for future benchmark development). Additionally, the study in this paper was conducted on a much larger scale (24 vs. 274 benchmarks), statistically highlighting the prevalent issues in existing benchmarks.

3 Guideline Design

3.1 The Lifecycle of Benchmark Development

Code-related benchmark development comprises five typical phases (Phase 1 - 5), as shown in Figure 1, explained in detail as follows.

Phase 1. Design. At the beginning of benchmark development, it is vital to identify the motivation, the *scope* and the *capabilities* required by the *application* scenario of interest. To achieve this

objective, one needs to carefully consider the application scenarios, making sure these scenarios align with real-world demands (Moriarty, 2011). Also, it is necessary to assess whether other benchmarks already exist that address similar tasks, and to identify any shortcomings they may possess (McIntosh et al., 2025; Malode, 2024). Furthermore, this new benchmark should be designed to evaluate specific LLMs’ capabilities; the crafted tasks are expected to reflect these capabilities (Hodak et al., 2023).

Phase 2. Construction. Benchmark Construction phase moves from design to execution. Typically, data is *collected* from public coding websites such as GitHub, LeetCode, and StackOverflow. This is followed by preprocessing, which includes filtering, cleaning (e.g., deduplication, denoising), and *curation* (e.g., aligning tests with corresponding code). The phase usually ends with a *validation* process, which can be manual or automated (McIntosh et al., 2025).

Phase 3. Evaluation. Once the benchmark is available, the next step is to apply it to LLMs, validating if it can effectively measure the intended LLM capabilities. Essential considering factors include *selecting a representative array of LLMs*, configuring *settings* like prompts and hyperparameters for consistency, choosing appropriate experimental *environments* to meet LLM requirements, and implementing thorough *logging* to ensure dependable and reproducible results.

Phase 4. Analysis. After evaluation, experimental results are analyzed, drawing conclusions on LLMs’ capabilities. This phase involves comparing each LLM’s *performance* to identify standout or underperforming models. Then, proper visual aids such as bar charts and tables can be used to *display* the experimental results, presenting clearer observation and deeper *inspiration*, such as the correlations between models, the correlations with re-

lated benchmarks, or performance in upper-/downstream tasks (Hendee and Wells, 1997). Indeed, a thorough analysis helps pinpoint areas for improvement and guides future LLM enhancements.

Phase 5. Release. The final phase is to make the benchmark open-accessible. This phase involves meticulously preparing all *materials* associated with the benchmark, ensuring they are ready for *open access* to foster widespread adoption and collaboration. Clear, comprehensive *documentation* is provided to guide users on effectively utilizing the benchmark. Additionally, all logged *experiment details* are made available, enhancing the reproducibility and transparency of the benchmark.

3.2 Study Design

Our study consists of four steps (Figure 2). All steps are explained as follows.

Step 1. Guideline Construction. To begin with, we *sketched the initial guidelines* for each phase in the benchmark development lifecycle (Section 3.1, Figure 1) by reviewing existing literature (Suppes et al., 1962; Zheng et al., 2023b; Schäfer et al., 2024; Reuel et al., 2024) and brainstorming. After that, we *refined the guidelines* through a series of interviews with various stakeholders, including model developers and benchmark builders, allowing for the addition, deletion, or modification of criteria based on expert feedback and practical insights. In order to allow more flexibility and increase the practicality of the guideline, we prioritized them with ★★★ (Highly important), ★★ (Important), and ★ (Optional). By doing so, the highly-recommended criteria are essential to comply with when setting up a new benchmark, while other criteria allow compromise, making the guideline simpler for benchmark developers to follow.

Step 2. Literature Profiling. This step begins by *collecting related benchmarks* according to their publication time, venue, and coding tasks, and then employing techniques like *snowballing* to ensure a comprehensive collection. This step leads to 274 code-related benchmarks for study. The detailed statistics can be found in the Appendix E. This step is followed by *profiling* each selected benchmark through a thorough review of *corresponding papers* and examination of *the released artifacts* or homepages associated with these benchmarks. The phase is completed by *reporting statistics* that highlight overall trends, pros, and cons identified during the profiling, providing a structured overview of existing benchmarks.

Step 3. Focused Case Study. After obtaining an overall impression of existing benchmarks, we *selected 30* ($= 5 * 6$) *representative benchmarks* from top-5 tasks, with top-5 highly-cited benchmarks plus the latest 1 benchmark (Appendix D). Each selected benchmark is then *analyzed against* HOW2BENCH, examining how well they meet the established criteria, studying their overall statistics, and identifying both exemplary and poor cases. Insights and references from existing literature are also incorporated to enrich the analysis, providing a deeper understanding of the benchmarks’ performance and areas for improvement.

Step 4. Human Study. The final step is a human study that evaluates the importance and practicality of HOW2BENCH. This involves *designing a questionnaire* by first initiating and iterating to gather diverse, logical insights, which is then *distributed* to a targeted audience. After collecting and filtering responses for quality, the data is *analyzed* to derive insights. See Appendix C for details.

4 Guideline – “HOW2BENCH”

The completed guideline HOW2BENCH with 55 criteria can be found in the Appendix. Each guideline contains: an actionable check criterion with necessary explanations; a priority indicator, divided into three levels in total, marking the importance of the checklist; and a checkbox for convenience.

4.1 Guideline for Benchmark Design

Explanation – For benchmark design, we listed four essential criteria, as shown in Figure 3. In particular, the guideline starts by recommending that benchmarks should initially assess if they are addressing a *significant gap* in existing research, ensuring the relevance and necessity of the benchmark. The *scope* of the benchmark is expected to be well-defined, clarifying the *capabilities* or characteristics being tested, how these relate to practical scenarios such as programming assistance or automated testing, and the relevance of these capabilities in real-world *applications*.

■ **Key Statistics** – According to our statistics among 270+ benchmarks, *apparent research bias* can be observed in terms of coding tasks, programming languages, and code granularities (Appendix B.1). For example, 36.13% (99/274) are code generation benchmarks, followed by program repair, with 9.85% (27/274). Also, more than 30 programming languages and 20 natural languages

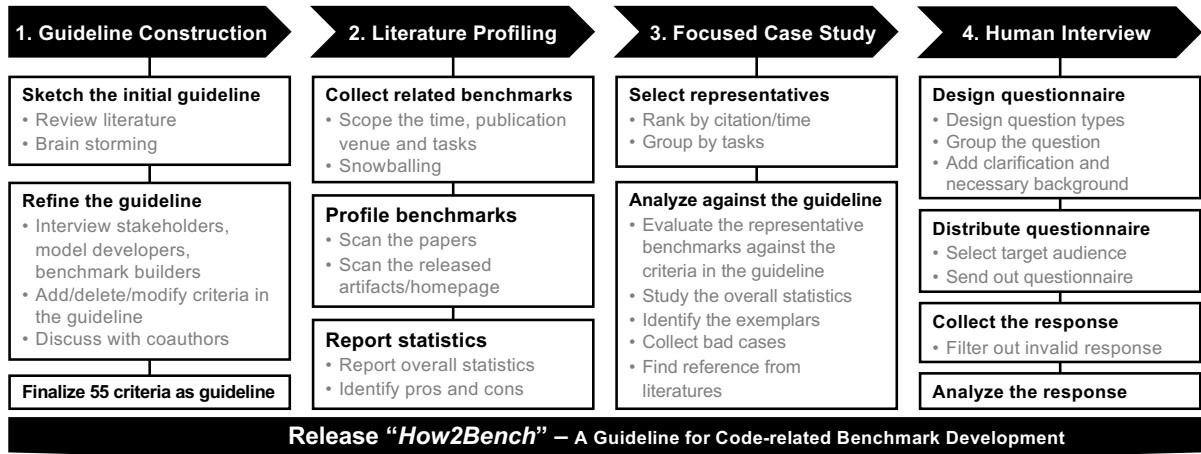


Figure 2: Workflow of study process

are assessed by only one benchmark.

Also, during the focused case study (listed in Appendix D), we identified that 10% benchmarks have not explicitly specified the capabilities (e.g., intention understanding, program synthesis) to be evaluated, and 30% *have not specified application scenarios* the benchmark targets.

Besides, we also identified a case in MBPP (Austin et al., 2021) where a case fell out of the target evaluation capabilities (Appendix B.2). Indeed, clearly defining the application scenarios/scopes/capabilities could help benchmark constructors establish precise goals for the design and development of the benchmark, ensuring accuracy in the evaluation.

Lastly, Figure 12 shows that 58% (158/274) code-related benchmarks involve Python, followed by 39% (107/274) involving Java. Yet, 31 programming languages are only covered by one benchmark, and less than five benchmarks cover other 19 programming languages. This observation consolidates the observation from previous works (Cao et al., 2024a; Hou et al., 2023) on a larger scale.

▲ **Severity** – Current benchmarks exhibit *an apparent imbalance* in coding tasks and programming languages dominated by code generation and Python, leaving research blank to be filled. Also, even highly cited benchmarks may have samples that do not fall into the examined capabilities.

4.2 Guideline for Construction

📖 **Explanation** – Figure 5 shows 19 criteria for benchmark construction. Essentially, for **data**

source, the key considerations include verifying the traceability and quality of the data source, addressing potential *data contamination* (Sainz et al., 2023), and ensuring that the *data sampling processes* are scientifically robust and rigorous. Also, for **data representativeness**, it also guides through specific checks to ensure the benchmark’s scope is strictly adhered to, such as making sure every data point falls within the targeted scope and that the data can cover all studied capabilities, domain knowledge, and application scenarios.

For data preprocessing and cleaning, it also stresses handling code-specific aspects, such as compilability and execution, along with cleaning and *manually reviewing* data for quality assurance. Output validation methods and evaluation metrics must be carefully designed and reviewed to ensure they effectively measure the benchmark’s goals. Lastly, it suggests considering additional evaluation perspectives, such as safety (Wei et al., 2024; Yuan et al., 2024a) checks, ensuring the code does not contain sensitive information.

■ **Key Statistics** – According to our statistics (Appendix B.3), the 270+ benchmarks exhibit *numerous irregularities* in their implementation, which could significantly threaten the reliability of the benchmarks. Surprisingly, **62% of benchmarks did not deduplicate** or did not mention. **Near 80% benchmarks did not consider or handle data contamination threats.** About 70% of the benchmarks did not go through any quality assurance checks such as manual checks and code execution. In particular, we summarized the *commonly-used data quality assurance metrics* and their frequency: manual check (22.6%), code execution

(2.2%), LLM check (1.5%), others (e.g. the number of stars or heuristic rules, 5.8%).

Also, since we focus on code-related benchmarks, which usually accompany test cases, *test coverage* also needs to be considered. As pointed out by a prior study (Liu et al., 2023a), inadequate test coverage can lead to inflated evaluation results. However, we observed that only 8.7% of benchmarks have considered test coverage when using test cases as oracles (Appendix B.3). It severely affects the reliability of findings on these benchmarks, potentially misleading future research and applications based on these flawed assessments.

▲ **Severity** – Most benchmarks display *severe loopholes* in data preparation and curation, i.e., only 31% of benchmarks went through a quality assurance check, and only 17% of them considered and handled data contamination threats.

4.3 Guideline for Evaluation

📖 **Explanation** – Guidelines for benchmark evaluation focus on the rigorousness and reliability of the evaluation. HOW2BENCH provides 12 criteria for benchmark evaluation, as shown in Figure 4. It mainly focuses on the comprehensive evaluation processes for benchmarks involving LLMs. For evaluation design, it stresses the importance of assessing a sufficient and *representative range of LLMs* to ensure the benchmark’s applicability across various model families and configurations, both open and closed-source. Figure 29 and Figure 30 show the distribution of numbers of LLMs studied and the most exercised LLMs.

Also, *prompting* has a direct impact on the quality of the LLMs’ output results (Wei et al., 2022; He et al., 2024a; Jin et al., 2024; Ye et al., 2023). As pointed out by a recent study, up to 40% performance gap could be observed in code translation when prompts vary (He et al., 2024b).

Additionally, *the experiment environment* is essential for reproducibility and transparency. Indeed, the hardware, software, and platform environments used during experiments might influence the outcomes (Ghosh, 2024). Furthermore, because of the nondeterministic nature of LLMs, experiments should be repeated, and randomization strategies should be used to mitigate the effects of randomness and parameter configuration biases. Lastly, *meticulously documented logs* of the experimental process is advised to facilitate transparency and reproducibility, detailing everything from param-

eter settings to the specific LLM pipelines such as vLLM (Kwon et al., 2023) used.

■ **Key Statistics** – Among the 274 benchmarks, 183 of them are evaluated over LLMs. According to our statistics (Figure 29), over 34% of the benchmarks were evaluated on fewer than 3 LLMs, with **11.48% benchmarks only evaluated on one LLM**. Such evaluation results can hardly be generalized to other LLMs. Furthermore, *more than half of the benchmarks studied fewer than 6 LLMs* ($51\% = (21 + 22 + 20 + 4 + 12 + 15)/183$).

For reference, we listed the top 10 most studied LLM families in Figure 30. Among them, the GPT series from OpenAI is the most extensively studied, accounting for 63% (116/183). Under the constraints of time and available resources, it is beneficial to evaluate more representative LLMs.

The prompt quality also significantly impacts the LLM evaluation (He et al., 2024b). According to a recent study, up to 40% of performance variation could be observed in the code translation task (He et al., 2024b). So, carefully designing a prompt needs consideration. However, **73.3%** representative benchmarks (Appendix D) do not validate whether the prompts they used are well-designed (Appendix B.4). Similarly, though 94.9% benchmarks were evaluated in a zero-shot manner, only 21.2% benchmarks were evaluated under few-shot, 8.8% under Chain-of-Thought and 2.6% under RAG (Appendix B.4). However, as shown in Figure 34, 73.3% representative benchmarks (Appendix D) do not validate whether the prompt they used is well-designed.

Regarding the evaluation process, our statistics exposed that **only 35.4% of benchmark evaluations have been repeated** (Appendix B.4). Also, regarding the transparency and matriculated documents, the observation is not optimistic – **Only 3.6% benchmarks provided their experiment environment. More than 50% of benchmarks did not provide reproducible instructions** such as prompts, examples for few-shot learning, or content for retrieval (Figure 39). **Less than half (42.7%) provide hyperparameters** such as temperature for reproduction.

▲ **Severity** – Over 60% of evaluations have not been repeated to eliminate the impact of randomness. **Less than 3.6% provided the complete and necessary information for reproducibility** such as prompts.

4.4 Guideline for Evaluation Analysis

☞ **Explanation** – The analysis of the experiment results is expected to be objective and comprehensive, hopefully providing insights or actionable advice. So, we listed 10 criteria for the evaluation analysis phase, as shown in Figure 6. Regarding **the perspectives of analysis**, inspired by classic measurement theory (Suppes et al., 1962), we suggest four essential perspectives, including **difficulty** (whether a benchmark is appropriately challenging for LLMs), **stability** (whether the results are consistent through repeated trials), **differentiability** (whether benchmarks can differentiate the strengths and weaknesses of various LLMs), and **inspiration** (e.g., the correlations between the upper-/downstream coding tasks and LLM scores).

Moreover, effective **presentation of results** using clear visual and textual descriptions could ensure the findings are understandable and actionable. The phase concludes with the suggestion to interpret and explain the results comprehensively, providing a basis for future enhancements.

■ **Key Statistics** – Because experimental analysis is relatively subjective and cannot be obtained through mechanical scanning, we focus on 30 representative focus benchmarks (Appendix D), covering the highest cited and latest benchmarks in top five tasks. Figure 37 shows an example from CruxEval (Gu et al., 2024) where the experimental scores can hardly be read from the figures.

Also, **explaining experiment results** is crucial for other practitioners to understand what the outcomes mean in the context of the research questions. According to our statistics (Appendix B.5), 70% benchmarks have detailed explanations and analyses of their evaluation results, while still **30% have not**. Indeed, an explanation contributes to the body of knowledge by making it possible to understand and compare results with previous studies, promoting transparency within the community.

▲ **Severity** – The analysis of experimental data and the clarity of data presentation may receive less attention and worth consideration. Even in papers cited 1k+ times like MBPP (Austin et al., 2021), there are instances of **unclear evaluation analysis and display**.

4.5 Guideline for Benchmark Release

☞ **Explanation** – Finally, releasing a benchmark for open access also needs careful consideration.

We offered 10 suggestions for this step, as shown in Figure 7, to highlight essential steps for public release preparation, emphasizing accessibility and ethical compliance. This includes setting an appropriate **license** to clarify usage rights, conducting a thorough review to **eliminate sensitive or harmful content** such as the API keys to access LLMs, the personal emails or toxic code comments (Miller et al., 2022) unless they are a part of the benchmark, and ensuring **transparency and reproducibility** by making all related materials openly available. **Detailed prompts** and clear descriptions of the experimental setup are advised to facilitate replication. Additionally, providing user manuals and evaluation interfaces is crucial for effective user engagement with the benchmark, enhancing its reliability and value for the research community.

■ **Key Statistics** – The final step involves the release of the benchmark. The fundamental requirement for releasing a benchmark is that it must be open-sourced. However, surprisingly, we observed that 5.1% of the benchmarks are only partially open-sourced (e.g., missing some subjects or tests), and **5.8% are not open-sourced at all** (e.g., links/web pages are no longer active). **19.3% have not properly set up the license**. Furthermore, prompts, which are necessary for reproducibility, are not disclosed in 52.6% of the benchmarks (Figure 39). Not to mention the lack of public information on experimental settings (Figure 32 and Figure 31) and experimental parameters (Figure 43). What is worse, 19.3% benchmarks do not setup licenses (Figure 44). The absence of licensing may lead to severe legal and ethical issues, potentially resulting in unauthorized use and distribution of proprietary technologies. Additionally, only 16.7% of the benchmarks make their logged experimental results publicly available (Appendix B.6).

▲ **Severity** – The release of existing benchmarks exhibits several issues. For example, over 10% of the benchmarks are either not open to public access or are only partially open-sourced. Only 47.4% of benchmarks are released with replicable prompts.

5 Human Study

To delve deeper into the integration of knowledge and action, we **surveyed 49 global researchers** in AI (42.6%) and SE (57.14%), as shown in Figure 50. Each participant had published at least one

research paper, and about *half had constructed code-related benchmarks*. See Appendix C.

First, **all participants agreed** that having a checklist for benchmark construction would contribute to the quality of the benchmark. 47/55 criteria in HOW2BENCH are deemed important by more 80% participants. Additionally, among the 21 participants who have constructed code-related benchmarks, *53 out of 55 criteria were deemed important by all benchmark developers*; only two criteria (criteria 3 and 4 in Section 4) were considered unimportant by a few individuals (3 and 2 participants, respectively). Additionally, we received two valuable suggestions that draw importance to recording *the time/monetary costs* of constructing the benchmark and conducting the experiments.

However, we also identified some *notable gaps in awareness*. First, regarding the *data preparation*, more than 15% of participants were not aware that the selection of data should consider the target scope of the evaluation set (i.e., the data must be representative), and 16% of participants were *unaware of the need for data denoising*. This oversight can significantly affect the validity and generalizability of experimental results, underscoring the importance of a comprehensive understanding of data handling for reliable research outcomes. Second, regarding *evaluation replicability and reliability*. Over 40% of participants believe that recording and publicizing the hardware and software environments, software versions, and libraries used in experiments is not important, with more than 20% still considering it unimportant despite already done so. This reveals *a significant lack of awareness* about the impact that experimental environments can have on *the reliability, reproducibility, and stability of evaluation results*. In fact, various studies have demonstrated that different experimental environments, parameters, and prompts can lead to substantial variations in outcomes (Xiao et al., 2024; Wang et al., 2019, 2023a).

6 Discussion

6.1 Trade-offs Between Benchmark Rigor and Development Efficiency

Despite all the above arguments, we admit that constructing rigorous benchmarks often entails a significant investment of time and human effort, which can lead to reduced efficiency in the development and evaluation process. Indeed, ensuring data quality, implementing thorough validation pro-

cedures, and designing comprehensive evaluation protocols requires non-trivial effort and may slow down the pace of research. It also echoed our human study. Through our human study, we found that researchers are often aware of the importance of several criteria (e.g., data denoise, and repeating the experiments), but did not implement them due to time constraints or other limitations.

However, this trade-off between rigor and efficiency is necessary to guarantee the reliability, reproducibility, and scientific value of benchmark results. While faster, less rigorous benchmarks might accelerate short-term experimentation, they risk producing misleading or non-generalizable findings that ultimately hinder long-term progress. Therefore, each checklist in HOW2BENCH is labeled with a **priority**: ★★★ (Highly important), ★★ (Important), and ★ (Optional). We hope the design of this indicator may better help benchmark developers balance efficiency with rigor.

6.2 Awareness and Action

It is also worth recapping the notable gaps in awareness of the importance of data preparation and reproducibility (Section 5). For example, 16% of participants were unaware of the need for data denoising; 40% of participants believe that recording and publicizing the hardware and software environments, software versions, and libraries used in experiments is not important.

However, without being aware of and addressing the quality of LLM benchmarks, the reliability and reproducibility of benchmark results remain compromised. This observation should be a call to action for the research community to strengthen education and awareness around best practices in benchmark development.

7 Conclusion

This paper proposes a rigorous guideline consisting of 55 checklists covering the benchmark development lifecycle. After investigating over 270 code-related benchmarks, we exposed their merits and limitations and provided suggestions for improving them. Finally, our human study reveals the neglect of details that may affect the benchmark’s reliability. In the long run, HOW2BENCH helps to improve the overall quality of benchmarks in the community due to the propagation among benchmarks.

Limitations

This paper has two primary limitations that offer avenues for future research. First, *the collection of code-related benchmarks may be incomplete*. To minimize this limitation, we covered papers published over the last decade, and conducted multiple rounds of snowballing. Ultimately, we collected 274 benchmarks, which is comparable to the number included in recent surveys (Hou et al., 2023; Schäfer et al., 2024) in the field. Second, *the study involved substantial manual analysis*, which could lead to oversight and discrepancies in the statistical results. To mitigate this issue, we ensured that each benchmark was double-checked by at least two authors and underwent multiple rounds of iteration. Third, *the guidelines may not cover all the details*. Constructing a code-related benchmark involves numerous details, and some criteria are task-specific. To overcome this limitation, we iteratively refined the guidelines, interviewed practitioners, and tried to cover the entire benchmark development process as thoroughly as possible. Last, *the human study participants may exhibit subjectivity*. To address this limitation, we endeavored to include a broad range of practitioners and seasoned researchers with experience in both AI and SE, aiming for the least biased results possible.

References

Yash Agarwal, Devansh Batra, and Ganesh Bagler. 2020. [Building hierarchically disentangled language models for text generation with named entities](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 26–38. International Committee on Computational Linguistics.

Rajas Agashe, Srinivasan Iyer, and Luke Zettlemoyer. 2019. [Juice: A large scale distantly supervised dataset for open domain context-based code generation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 5435–5445. Association for Computational Linguistics.

Lakshya A. Agrawal, Aditya Kanade, Navin Goyal, Shuvendu K. Lahiri, and Sriram K. Rajamani. 2023. [Monitor-guided decoding of code lms with static analysis of repository context](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Wasi Uddin Ahmad, Md Golam Rahman Tushar, Saikat Chakraborty, and Kai-Wei Chang. 2023. [AVATAR: A parallel corpus for java-python program translation](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 2268–2281. Association for Computational Linguistics.

Miltiadis Allamanis, Henry Jackson-Flux, and Marc Brockschmidt. 2021. [Self-supervised bug detection and repair](#). In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 27865–27876.

Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav. 2019. [code2seq: Generating sequences from structured representations of code](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.

Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Sujan Kumar Gonugondla, Hantian Ding, Varun Kumar, Nathan Fulton, Arash Farahani, Siddhartha Jain, Robert Giaquinto, Haifeng Qian, Murali Krishna Ramanathan, Ramesh Nallapati, Baishakhi Ray, Parminder Bhatia, Sudipta Sengupta, Dan Roth, and Bing Xiang. 2022. [Multi-lingual evaluation of code generation models](#).

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Hannah McLean Babe, Sydney Nguyen, Yangtian Zi, Arjun Guha, Molly Q. Feldman, and Carolyn Jane Anderson. 2024. [Studenteval: A benchmark of student-written prompts for large language models of code](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 8452–8474. Association for Computational Linguistics.

Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Vageesh D. C., Arun Iyer, Suresh Parthasarathy, Sriram K. Rajamani, Balasubramanyan Ashok, and Shashank Shet. 2024. [Codeplan: Repository-level coding using llms and planning](#). *Proc. ACM Softw. Eng.*, 1(FSE):675–698.

- Antonio Valerio Miceli Barone and Rico Sennrich. 2017. [A parallel corpus of python functions and documentation strings for automated code documentation and code generation](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017, Volume 2: Short Papers*, pages 314–319. Asian Federation of Natural Language Processing.
- Earl T Barr, Mark Harman, Phil McMinn, Muzammil Shahbaz, and Shin Yoo. 2014. The oracle problem in software testing: A survey. *IEEE transactions on software engineering*, 41(5):507–525.
- Berkay Berabi, Jingxuan He, Veselin Raychev, and Martin T. Vechev. 2021. [Tfix: Learning to fix coding errors with a text-to-text transformer](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 780–791. PMLR.
- Egor Bogomolov, Aleksandra Eliseeva, Timur Galimzyanov, Evgeniy Glukhov, Anton Shapkin, Maria Tigina, Yaroslav Golubev, Alexander Kovrigin, Arie van Deursen, Maliheh Izadi, and Timofey Bryksin. 2024. [Long code arena: a set of benchmarks for long-context code models](#). *CoRR*, abs/2406.11612.
- Jialun Cao, Zhiyong Chen, Jiarong Wu, Shing-Chi Cheung, and Chang Xu. 2024a. [Can AI beat undergraduates in entry-level java assignments? benchmarking large language models on javabench](#). *CoRR*, abs/2406.12902.
- Jialun Cao, Wuqi Zhang, and Shing-Chi Cheung. 2024b. Concerned with data contamination? assessing countermeasures in code language model. *arXiv preprint arXiv:2403.16898*.
- Ruisheng Cao, Fangyu Lei, Haoyuan Wu, Jixuan Chen, Yeqiao Fu, Hongcheng Gao, Xinzhuang Xiong, Hanchong Zhang, Yuchen Mao, Wenjing Hu, Tianbao Xie, Hongshen Xu, Danyang Zhang, Sida Wang, Ruoxi Sun, Pengcheng Yin, Caiming Xiong, Ansong Ni, Qian Liu, Victor Zhong, Lu Chen, Kai Yu, and Tao Yu. 2024c. [Spider2-v: How far are multimodal agents from automating data science and engineering workflows?](#) *CoRR*, abs/2407.10956.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. 2022. Multipl-e: A scalable and extensible approach to benchmarking neural code generation. *arXiv preprint arXiv:2208.08227*.
- Saikat Chakraborty, Rahul Krishna, Yangruibo Ding, and Baishakhi Ray. 2022. [Deep learning based vulnerability detection: Are we there yet?](#) *IEEE Trans. Software Eng.*, 48(9):3280–3296.
- Shubham Chandel, Colin B Clement, Guillermo Serrato, and Neel Sundaresan. 2022. Training and evaluating a jupyter notebook data science assistant. *arXiv preprint arXiv:2201.12901*.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebguss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021a. [Evaluating large language models trained on code](#).
- Xinyun Chen, Linyuan Gong, Alvin Cheung, and Dawn Song. 2021b. [Plotcoder: Hierarchical decoding for synthesizing visualization code in programmatic context](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2169–2181. Association for Computational Linguistics.
- Yizheng Chen, Zhoujie Ding, Lamya Alowain, Xinyun Chen, and David A. Wagner. 2023. [Diversevul: A new vulnerable source code dataset for deep learning based vulnerability detection](#). In *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses, RAID 2023, Hong Kong, China, October 16-18, 2023*, pages 654–668. ACM.
- Jianbo Dai, Jianqiao Lu, Yunlong Feng, Rongju Ruan, Ming Cheng, Haochen Tan, and Zhijiang Guo. 2024. [MHPP: exploring the capabilities and limitations of language models beyond basic code generation](#). *CoRR*, abs/2405.11430.
- Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek, Oleksandr Polozov, Huan Sun, and Matthew Richardson. 2021. [Structure-grounded pretraining for text-to-sql](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 1337–1350. Association for Computational Linguistics.
- Yangruibo Ding, Yanjun Fu, Omniyyah Ibrahim, Chawin Sitawarin, Xinyun Chen, Basel Alomair,

920	David A. Wagner, Baishakhi Ray, and Yizheng Chen.	pages 351–360. Association for Computational Lin-	977
921	2024a. Vulnerability detection with code language	guistics.	978
922	models: How far are we? <i>CoRR</i> , abs/2403.18624.		
923	Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Han-	Karl Fogel. 2005. <i>Producing open source software:</i>	979
924	tian Ding, Ming Tan, Nihal Jain, Murali Krishna Ra-	<i>How to run a successful free software project.</i> "	980
925	manathan, Ramesh Nallapati, Parminder Bhatia, Dan	O'Reilly Media, Inc."	981
926	Roth, and Bing Xiang. 2023. Crosscodeeval: A di-		
927	verse and multilingual benchmark for cross-file code	Graham Francis and Jacky Holloway. 2007. What have	982
928	completion. In <i>Advances in Neural Information Pro-</i>	we learned? themes from the literature on best-	983
929	<i>cessing Systems 36: Annual Conference on Neural</i>	practice benchmarking. <i>International Journal of</i>	984
930	<i>Information Processing Systems 2023, NeurIPS 2023,</i>	<i>Management Reviews</i> , 9(3):171–189.	985
931	<i>New Orleans, LA, USA, December 10 - 16, 2023.</i>		
932	Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad,	Lingyue Fu, Huacan Chai, Shuang Luo, Kounianhua Du,	986
933	Murali Krishna Ramanathan, Ramesh Nallapati, Par-	Weiming Zhang, Longteng Fan, Jiayi Lei, Renting	987
934	minder Bhatia, Dan Roth, and Bing Xiang. 2024b.	Rui, Jianghao Lin, Yuchen Fang, Yifan Liu, Jingkuan	988
935	Cocomic: Code completion by jointly modeling in-	Wang, Siyuan Qi, Kangning Zhang, Weinan Zhang,	989
936	file and cross-file context. In <i>Proceedings of the</i>	and Yong Yu. 2023. Codeapex: A bilingual pro-	990
937	<i>2024 Joint International Conference on Computa-</i>	gramming evaluation benchmark for large language	991
938	<i>tional Linguistics, Language Resources and Evalua-</i>	models. <i>CoRR</i> , abs/2309.01940.	992
939	<i>tion, LREC/COLING 2024, 20-25 May, 2024, Torino,</i>		
940	<i>Italy</i> , pages 3433–3445. ELRA and ICCL.	YanJun Fu, Ethan Baker, and Yizheng Chen. 2024. Con-	993
		strained decoding for secure code generation. <i>CoRR</i> ,	994
		abs/2405.00218.	995
941	Shihan Dou, Haoxiang Jia, Shenxi Wu, Huiyuan Zheng,	Yujian Gan, Xinyun Chen, Qiuping Huang, and	996
942	Weikang Zhou, Muling Wu, Mingxu Chai, Jessica	Matthew Purver. 2022. Measuring and improving	997
943	Fan, Caishuang Huang, Yunbo Tao, et al. 2024.	compositional generalization in text-to-sql via com-	998
944	What's wrong with your code generated by large	ponent alignment. In <i>Findings of the Association</i>	999
945	language models? an extensive study. <i>arXiv preprint</i>	<i>for Computational Linguistics: NAACL 2022, Seattle,</i>	1000
946	<i>arXiv:2407.06153.</i>	<i>WA, United States, July 10-15, 2022</i> , pages 831–843.	1001
		Association for Computational Linguistics.	1002
947	Mingzhe Du, Anh Tuan Luu, Bin Ji, Qian Liu, and	Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew	1003
948	See-Kiong Ng. 2024. Mercury: A code efficiency	Purver, John R. Woodward, Jinxia Xie, and Peng-	1004
949	benchmark for code large language models. In <i>The</i>	sheng Huang. 2021a. Towards robustness of text-to-	1005
950	<i>Thirty-eight Conference on Neural Information Pro-</i>	sql models against synonym substitution. In <i>Proceed-</i>	1006
951	<i>cessing Systems Datasets and Benchmarks Track.</i>	<i>ings of the 59th Annual Meeting of the Association for</i>	1007
952	Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang,	<i>Computational Linguistics and the 11th International</i>	1008
953	Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng	<i>Joint Conference on Natural Language Processing,</i>	1009
954	Sha, Xin Peng, and Yiling Lou. 2023a. Classe-	<i>ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual</i>	1010
955	eval: A manually-crafted benchmark for evaluat-	<i>Event, August 1-6, 2021</i> , pages 2505–2515. Associa-	1011
956	ing llms on class-level code generation. <i>Preprint,</i>	<i>tion for Computational Linguistics.</i>	1012
957	<i>arXiv:2308.01861.</i>		
958	Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang,	Yujian Gan, Xinyun Chen, and Matthew Purver. 2021b.	1013
959	Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha,	Exploring underexplored limitations of cross-domain	1014
960	Xin Peng, and Yiling Lou. 2023b. Classeval: A	text-to-sql generalization. In <i>Proceedings of the 2021</i>	1015
961	manually-crafted benchmark for evaluating llms on	<i>Conference on Empirical Methods in Natural Lan-</i>	1016
962	class-level code generation. <i>CoRR</i> , abs/2308.01861.	<i>guage Processing, EMNLP 2021, Virtual Event /</i>	1017
		<i>Punta Cana, Dominican Republic, 7-11 November,</i>	1018
963	Aleksandra Eliseeva, Yaroslav Sokolov, Egor Bogo-	<i>2021</i> , pages 8926–8931. Association for Computa-	1019
964	molov, Yaroslav Golubev, Danny Dig, and Timofey	<i>tional Linguistics.</i>	1020
965	Bryksin. 2023. From commit message generation		
966	to history-aware commit message completion. In	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon,	1021
967	<i>38th IEEE/ACM International Conference on Auto-</i>	Pengfei Liu, Yiming Yang, Jamie Callan, and Gra-	1022
968	<i>mated Software Engineering, ASE 2023, Luxembourg,</i>	ham Neubig. 2023a. PAL: program-aided language	1023
969	<i>September 11-15, 2023</i> , pages 723–735. IEEE.	models. In <i>International Conference on Machine</i>	1024
		<i>Learning, ICML 2023, 23-29 July 2023, Honolulu,</i>	1025
970	Catherine Finegan-Dollak, Jonathan K. Kummerfeld,	<i>Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine</i>	1026
971	Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui	<i>Learning Research</i> , pages 10764–10799. PMLR.	1027
972	Zhang, and Dragomir R. Radev. 2018. Improving		
973	text-to-sql evaluation methodology. In <i>Proceedings</i>	Zeyu Gao, Hao Wang, Yuchen Zhou, Wenyu Zhu, and	1028
974	<i>of the 56th Annual Meeting of the Association for</i>	Chao Zhang. 2023b. How far have we gone in vulner-	1029
975	<i>Computational Linguistics, ACL 2018, Melbourne,</i>	ability detection using large language models. <i>CoRR</i> ,	1030
976	<i>Australia, July 15-20, 2018, Volume 1: Long Papers,</i>	abs/2311.12420.	1031

1143	ICPC 2018, Gothenburg, Sweden, May 27-28, 2018,	1199
1144	pages 200–210. ACM.	1200
1145	Xing Hu, Ge Li, Xin Xia, David Lo, Shuai Lu, and	1201
1146	Zhi Jin. 2018b. Summarizing source code with trans-	1202
1147	ferred API knowledge . In <i>Proceedings of the Twenty-</i>	1203
1148	<i>Seventh International Joint Conference on Artificial</i>	1204
1149	<i>Intelligence, IJCAI 2018, July 13-19, 2018, Stock-</i>	1205
1150	<i>holm, Sweden</i> , pages 2269–2275. ijcai.org.	
1151	Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli	1206
1152	Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing	1207
1153	Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li,	1208
1154	Kun Kuang, Yang Yang, Hongxia Yang, and Fei Wu.	1209
1155	2024. Infiagent-dabench: Evaluating agents on data	1210
1156	analysis tasks . In <i>Forty-first International Confer-</i>	1211
1157	<i>ence on Machine Learning, ICML 2024, Vienna, Aus-</i>	1212
1158	<i>tria, July 21-27, 2024</i> . OpenReview.net.	
1159	Yang Hu, Umair Z. Ahmed, Sergey Mechtaev, Ben	1213
1160	Leong, and Abhik Roychoudhury. 2019. Re-	1214
1161	factoring based program repair applied to program-	1215
1162	ming assignments . In <i>2019 34th IEEE/ACM Interna-</i>	1216
1163	<i>tional Conference on Automated Software Engineer-</i>	1217
1164	<i>ing (ASE)</i> , pages 388–398.	1218
1165	Dong HUANG, Yuhao QING, Weiyi Shang, Heming	1219
1166	Cui, and Jie Zhang. 2024. Effibench: Benchmarking	1220
1167	the efficiency of automatically generated code . In	1221
1168	<i>The Thirty-eight Conference on Neural Information</i>	1222
1169	<i>Processing Systems Datasets and Benchmarks Track</i> .	1223
1170	Junjie Huang, Duyu Tang, Linjun Shou, Ming Gong,	1224
1171	Ke Xu, Daxin Jiang, Ming Zhou, and Nan Duan.	1225
1172	2021. CoSQA: 20,000+ web queries for code search	1226
1173	and question answering . In <i>Proceedings of the 59th</i>	1227
1174	<i>Annual Meeting of the Association for Computational</i>	1228
1175	<i>Linguistics and the 11th International Joint Confer-</i>	1229
1176	<i>ence on Natural Language Processing (Volume 1:</i>	1230
1177	<i>Long Papers)</i> , pages 5690–5700, Online. Association	
1178	for Computational Linguistics.	
1179	Junjie Huang, Chenglong Wang, Jipeng Zhang, Cong	1231
1180	Yan, Haotian Cui, Jeevana Priya Inala, Colin B.	1232
1181	Clement, Nan Duan, and Jianfeng Gao. 2022.	1233
1182	Execution-based evaluation for data science code	1234
1183	generation models . <i>CoRR</i> , abs/2211.09374.	1235
1184	Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong,	1236
1185	Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong Shen,	1237
1186	Chen Lin, Nan Duan, and Weizhu Chen. 2024.	
1187	Competition-level problems are effective LLM eval-	1238
1188	uators . In <i>Findings of the Association for Computa-</i>	1239
1189	<i>tional Linguistics, ACL 2024, Bangkok, Thailand and</i>	1240
1190	<i>virtual meeting, August 11-16, 2024</i> , pages 13526–	1241
1191	13544. Association for Computational Linguistics.	1242
1192	Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis	1243
1193	Allamanis, and Marc Brockschmidt. 2019. Code-	
1194	searchnet challenge: Evaluating the state of semantic	1244
1195	code search . <i>CoRR</i> , abs/1909.09436.	1245
1196	Shima Imani, Liang Du, and Harsh Shrivastava. 2023.	1246
1197	Mathprompter: Mathematical reasoning using large	1247
1198	language models . <i>Preprint</i> , arXiv:2303.05398.	1248
	Marko Ivanković, Goran Petrović, René Just, and Gor-	1249
	don Fraser. 2019. Code coverage at google . In	1250
	<i>Proceedings of the 2019 27th ACM Joint Meeting</i>	1251
	<i>on European Software Engineering Conference and</i>	
	<i>Symposium on the Foundations of Software Engineer-</i>	1252
	<i>ing, ESEC/FSE 2019</i> , page 955–963, New York, NY,	1253
	USA. Association for Computing Machinery.	1254
	Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and	1255
	Luke Zettlemoyer. 2016. Summarizing source code	1256
	using a neural attention model . In <i>Proceedings of the</i>	
	<i>54th Annual Meeting of the Association for Compu-</i>	
	<i>tational Linguistics, ACL 2016, August 7-12, 2016,</i>	
	<i>Berlin, Germany, Volume 1: Long Papers</i> . The Asso-	
	ciation for Computer Linguistics.	
	Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and	
	Luke Zettlemoyer. 2018. Mapping language to code	
	in programmatic context . In <i>Proceedings of the 2018</i>	
	<i>Conference on Empirical Methods in Natural Lan-</i>	
	<i>guage Processing</i> , pages 1643–1652, Brussels, Bel-	
	gium. Association for Computational Linguistics.	
	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia	
	Yan, Tianjun Zhang, Sida Wang, Armando Solar-	
	Lezama, Koushik Sen, and Ion Stoica. 2024. Live-	
	codebench: Holistic and contamination free eval-	
	uation of large language models for code . <i>CoRR</i> ,	
	abs/2403.07974.	
	Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan.	
	2023. Impact of code language models on automated	
	program repair . In <i>45th IEEE/ACM International</i>	
	<i>Conference on Software Engineering, ICSE 2023,</i>	
	<i>Melbourne, Australia, May 14-20, 2023</i> , pages 1430–	
	1442. IEEE.	
	Mingsheng Jiao, Tingrui Yu, Xuan Li, Guanjie Qiu,	
	Xiaodong Gu, and Beijun Shen. 2023. On the eval-	
	uation of neural code translation: Taxonomy and	
	benchmark . In <i>38th IEEE/ACM International Con-</i>	
	<i>ference on Automated Software Engineering, ASE</i>	
	<i>2023, Luxembourg, September 11-15, 2023</i> , pages	
	1529–1541. IEEE.	
	Carlos E Jimenez, John Yang, Alexander Wettig,	
	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R	
	Narasimhan. 2024. SWE-bench: Can language mod-	
	els resolve real-world github issues? In <i>The Twelfth</i>	
	<i>International Conference on Learning Representa-</i>	
	<i>tions</i> .	
	Matthew Jin, Syed Shahriar, Michele Tufano, Xin	
	Shi, Shuai Lu, Neel Sundaresan, and Alexey Svy-	
	atkovskiy. 2023. Inferfix: End-to-end program repair	
	with llms . In <i>Proceedings of the 31st ACM Joint Eu-</i>	
	<i>ropean Software Engineering Conference and Sym-</i>	
	<i>posium on the Foundations of Software Engineering,</i>	
	<i>ESEC/FSE 2023, San Francisco, CA, USA, Decem-</i>	
	<i>ber 3-9, 2023</i> , pages 1646–1656. ACM.	
	Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao,	
	Wenyue Hua, Yanda Meng, Yongfeng Zhang, and	
	Mengnan Du. 2024. The impact of reasoning step	
	length on large language models. <i>arXiv preprint</i>	
	<i>arXiv:2401.04925</i> .	

- René Just, Darioush Jalali, and Michael D. Ernst. 2014. [Defects4j: a database of existing faults to enable controlled testing studies for java programs](#). In *International Symposium on Software Testing and Analysis, ISSTA '14, San Jose, CA, USA - July 21 - 26, 2014*, pages 437–440. ACM.
- Mohammad Abdullah Matin Khan, M. Saiful Bari, Xuan Do Long, Weishi Wang, Md. Rizwan Parvez, and Shafiq Joty. 2024. [Xcodeeval: An execution-based large scale multilingual multitask benchmark for code understanding, generation, translation and retrieval](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 6766–6805. Association for Computational Linguistics.
- Rahul Kumar, Amar Raja Dibbu, Shrutendra Harsola, Vignesh Subrahmaniam, and Ashutosh Modi. 2024. [Booksql: A large scale text-to-sql dataset for accounting domain](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico, June 16-21, 2024, pages 497–516. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Beck LaBash, August Rosedale, Alex Reents, Lucas Negritto, and Colin Wiel. 2024. [RES-Q: evaluating code-editing large language model systems at the repository scale](#). *CoRR*, abs/2406.16801.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-Tau Yih, Daniel Fried, Sida I. Wang, and Tao Yu. 2023. [DS-1000: A natural and reliable benchmark for data science code generation](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 18319–18345. PMLR.
- Claire Le Goues, Neal J. Holtschulte, Edward K. Smith, Yuriy Brun, Premkumar T. Devanbu, Stephanie Forrest, and Westley Weimer. 2015. [The manybugs and introclass benchmarks for automated repair of C programs](#). *IEEE Trans. Software Eng.*, 41(12):1236–1256.
- Alexander LeClair, Siyuan Jiang, and Collin McMillan. 2019. [A neural model for generating natural language summaries of program subroutines](#). In *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*, pages 795–806. IEEE / ACM.
- Changyoon Lee, Yeon Seonwoo, and Alice Oh. 2022. [CSIQA: 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, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 2026–2040. Association for Computational Linguistics.
- Chia-Hsuan Lee, Oleksandr Polozov, and Matthew Richardson. 2021. [Kaggledbqa: Realistic evaluation of text-to-sql parsers](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2261–2273. Association for Computational Linguistics.
- Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2023. [EHRSQL: A practical text-to-sql benchmark for electronic health records](#). *CoRR*, abs/2301.07695.
- Jia Li, Ge Li, Yunfei Zhao, Yongmin Li, Huanyu Liu, Hao Zhu, Lecheng Wang, Kaibo Liu, Zheng Fang, Lanshen Wang, Jiazheng Ding, Xuanming Zhang, Yuqi Zhu, Yihong Dong, Zhi Jin, Binhua Li, Fei Huang, and Yongbin Li. 2024a. [DevEval: A Manually-Annotated Code Generation Benchmark Aligned with Real-World Code Repositories](#). *arXiv preprint*. ArXiv:2405.19856 [cs].
- Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin Chen-Chuan Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2023a. [Can LLM already serve as A database interface? A big bench for large-scale database grounded text-to-sqls](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Kaixin Li, Qisheng Hu, James Xu Zhao, Hui Chen, Yuxi Xie, Tiedong Liu, Michael Shieh, and Junxian He. 2024b. [Instructcoder: Instruction tuning large language models for code editing](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, ACL 2024 - Student Research Workshop, Bangkok, Thailand, August 11-16, 2024*, pages 50–70. Association for Computational Linguistics.
- Kaixin Li, Yuchen Tian, Qisheng Hu, Ziyang Luo, and Jing Ma. 2024c. [Mmcode: Evaluating multi-modal code large language models with visually rich programming problems](#). *CoRR*, abs/2404.09486.
- Linyi Li, Shijie Geng, Zhenwen Li, Yibo He, Hao Yu, Ziyue Hua, Guanghan Ning, Siwei Wang, Tao Xie, and Hongxia Yang. 2024d. [Infibench: Evaluating](#)

1373	the question-answering capabilities of code large language models. <i>Preprint</i> , arXiv:2404.07940.	1431
1374		1432
1375	Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. 2023b. TACO: topics in algorithmic code generation dataset. <i>CoRR</i> , abs/2312.14852.	1433
1376		1434
1377		1435
1378		
1379	Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. Competition-level code generation with alpha-code. <i>Science</i> , 378(6624):1092–1097.	1436
1380		1437
1381		1438
1382		1439
1383		1440
1384		
1385		1441
1386		1442
1387		1443
1388		1444
1389		1445
1390	Zehan Li, Jianfei Zhang, Chuantao Yin, Yuanxin Ouyang, and Wenge Rong. 2024e. Procqa: A large-scale community-based programming question answering dataset for code search. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy</i> , pages 13057–13067. ELRA and ICCL.	1446
1391		1447
1392		1448
1393		1449
1394		1450
1395		1451
1396		1452
1397		
1398	Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Yawei Zhu, Zhaoxuan Chen, Sujuan Wang, and Jialai Wang. 2018a. Sysevr: A framework for using deep learning to detect software vulnerabilities. <i>CoRR</i> , abs/1807.06756.	1453
1399		1454
1400		1455
1401		1456
1402		
1403	Zhen Li, Deqing Zou, Shouhuai Xu, Xinyu Ou, Hai Jin, Sujuan Wang, Zhijun Deng, and Yuyi Zhong. 2018b. Vuldeepecker: A deep learning-based system for vulnerability detection. In <i>25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018</i> . The Internet Society.	1457
1404		1458
1405		1459
1406		1460
1407		1461
1408		1462
1409		
1410	Dianshu Liao, Shidong Pan, Xiaoyu Sun, Xiaoxue Ren, Qing Huang, Zhenchang Xing, Huan Jin, and Qinying Li. 2024. A 3-codgen: A repository-level code generation framework for code reuse with local-aware, global-aware, and third-party-library-aware. <i>IEEE Transactions on Software Engineering</i> .	1463
1411		1464
1412		1465
1413		1466
1414		1467
1415		1468
1416	Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. 2017. Quixbugs: a multi-lingual program repair benchmark set based on the quixey challenge. In <i>Proceedings Companion of the 2017 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity, SPLASH 2017, Vancouver, BC, Canada, October 23 - 27, 2017</i> , pages 55–56. ACM.	1469
1417		1470
1418		
1419		1471
1420		1472
1421		1473
1422		1474
1423		1475
1424		1476
1425	Guanjun Lin, Wei Xiao, Jun Zhang, and Yang Xiang. 2019. Deep learning-based vulnerable function detection: A benchmark. In <i>Information and Communications Security - 21st International Conference, ICICS 2019, Beijing, China, December 15-17, 2019, Revised Selected Papers</i> , volume 11999 of <i>Lecture Notes in Computer Science</i> , pages 219–232. Springer.	1477
1426		1478
1427		1479
1428		1480
1429		
1430		1481
		1482
		1483
		1484
		1485
		1486
	Guanjun Lin, Jun Zhang, Wei Luo, Lei Pan, Olivier Y. de Vel, Paul Montague, and Yang Xiang. 2021. Software vulnerability discovery via learning multi-domain knowledge bases. <i>IEEE Trans. Dependable Secur. Comput.</i> , 18(5):2469–2485.	1431
		1432
		1433
		1434
		1435
	Guanjun Lin, Jun Zhang, Wei Luo, Lei Pan, Yang Xiang, Olivier Y. de Vel, and Paul Montague. 2018. Cross-project transfer representation learning for vulnerable function discovery. <i>IEEE Trans. Ind. Informatics</i> , 14(7):3289–3297.	1436
		1437
		1438
		1439
		1440
	Xiang Ling, Lingfei Wu, Saizhuo Wang, Gaoning Pan, Tengfei Ma, Fangli Xu, Alex X. Liu, Chunming Wu, and Shouling Ji. 2021. Deep graph matching and searching for semantic code retrieval. <i>ACM Trans. Knowl. Discov. Data</i> , 15(5):88:1–88:21.	1441
		1442
		1443
		1444
		1445
	Chenxiao Liu and Xiaojun Wan. 2021. Codeqa: A question answering dataset for source code comprehension. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021</i> , pages 2618–2632. Association for Computational Linguistics.	1446
		1447
		1448
		1449
		1450
		1451
		1452
	Jiawei Liu, Jia Le Tian, Vijay Daita, Yuxiang Wei, Yifeng Ding, Yuhang Katherine Wang, Jun Yang, and Lingming Zhang. 2024a. Repoqa: Evaluating long context code understanding. <i>CoRR</i> , abs/2406.06025.	1453
		1454
		1455
		1456
	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023a. Is your code generated by chat-GPT really correct? rigorous evaluation of large language models for code generation. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1457
		1458
		1459
		1460
		1461
		1462
	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023b. Is your code generated by chat-gpt really correct? rigorous evaluation of large language models for code generation. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
	Shangqing Liu, Yu Chen, Xiaofei Xie, Jing Kai Siow, and Yang Liu. 2021. Retrieval-augmented generation for code summarization via hybrid GNN. In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	1471
		1472
		1473
		1474
		1475
		1476
	Shangqing Liu, Cuiyun Gao, Sen Chen, Lun Yiu Nie, and Yang Liu. 2022. ATOM: commit message generation based on abstract syntax tree and hybrid ranking. <i>IEEE Trans. Software Eng.</i> , 48(5):1800–1817.	1477
		1478
		1479
		1480
	Tianyang Liu, Canwen Xu, and Julian J. McAuley. 2024b. Repobench: Benchmarking repository-level code auto-completion systems. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1481
		1482
		1483
		1484
		1485
		1486

1487	Yu Liu, Lang Gao, Mingxin Yang, Yu Xie, Ping Chen,	deep similarity learning-based type inference for	1544
1488	Xiaojin Zhang, and Wei Chen. 2024c. Vuldetect-	python . In <i>44th IEEE/ACM 44th International Con-</i>	1545
1489	bench: Evaluating the deep capability of vulnera-	ference on Software Engineering, ICSE 2022, Pitts-	1546
1490	bility detection with large language models . <i>CoRR</i> ,	burgh, PA, USA, May 25-27, 2022, pages 2241–2252.	1547
1491	abs/2406.07595.	ACM.	1548
1492	Yuliang Liu, Xiangru Tang, Zefan Cai, Junjie Lu,	John P Moriarty. 2011. A theory of benchmarking.	1549
1493	Yichi Zhang, Yanjun Shao, Zexuan Deng, Helan Hu,	<i>Benchmarking: An International Journal</i> , 18(4):588–	1550
1494	Zengxian Yang, Kaikai An, Ruijun Huang, Shuzheng	611.	1551
1495	Si, Sheng Chen, Haozhe Zhao, Zhengliang Li, Liang	Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016.	1552
1496	Chen, Yiming Zong, Yan Wang, Tianyu Liu, Zhi-	Convolutional neural networks over tree structures	1553
1497	wei Jiang, Baobao Chang, Yujia Qin, Wangchunshu	for programming language processing . In <i>Proceed-</i>	1554
1498	Zhou, Yilun Zhao, Arman Cohan, and Mark Ger-	ings of the Thirtieth AAAI Conference on Artificial	1555
1499	stein. 2023c. ML-bench: Large language models	Intelligence, February 12-17, 2016, Phoenix, Ari-	1556
1500	leverage open-source libraries for machine learning	zona, USA, pages 1287–1293. AAAI Press.	1557
1501	tasks . <i>CoRR</i> , abs/2311.09835.		
1502	Zhongxin Liu, Xin Xia, Ahmed E. Hassan, David Lo,	Hussein Mozannar, Valerie Chen, Mohammed Alsobay,	1558
1503	Zhenchang Xing, and Xinyu Wang. 2018. Neural-	Subhro Das, Sebastian Zhao, Dennis Wei, Manish	1559
1504	machine-translation-based commit message gener-	Nagireddy, Prasanna Sattigeri, Ameet Talwalkar, and	1560
1505	ation: how far are we? In <i>Proceedings of the</i>	David A. Sontag. 2024. The realhumaneval: Eval-	1561
1506	<i>33rd ACM/IEEE International Conference on Auto-</i>	uating large language models’ abilities to support	1562
1507	<i>ated Software Engineering, ASE 2018, Montpellier,</i>	programmers . <i>CoRR</i> , abs/2404.02806.	1563
1508	<i>France, September 3-7, 2018, pages 373–384. ACM.</i>		
1509	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey	Niklas Muennighoff, Qian Liu, Armel Randy Ze-	1564
1510	Svyatkovskiy, Ambrosio Blanco, Colin B. Clement,	baze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo,	1565
1511	Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Li-	Swayam Singh, Xiangru Tang, Leandro von Werra,	1566
1512	dong Zhou, Linjun Shou, Long Zhou, Michele Tu-	and Shayne Longpre. 2024. Octopack: Instruction	1567
1513	fano, Ming Gong, Ming Zhou, Nan Duan, Neel Sun-	tuning code large language models . In <i>The Twelfth</i>	1568
1514	daresan, Shao Kun Deng, Shengyu Fu, and Shujie	<i>International Conference on Learning Representa-</i>	1569
1515	Liu. 2021. Codexglue: A machine learning bench-	tions, ICLR 2024, Vienna, Austria, May 7-11, 2024.	1570
1516	mark dataset for code understanding and generation .	OpenReview.net.	1571
1517	In <i>Proceedings of the Neural Information Process-</i>	Hoan Anh Nguyen, Tien N. Nguyen, Danny Dig,	1572
1518	<i>ing Systems Track on Datasets and Benchmarks 1,</i>	Son Nguyen, Hieu Tran, and Michael Hilton. 2019.	1573
1519	<i>NeurIPS Datasets and Benchmarks 2021, December</i>	Graph-based mining of in-the-wild, fine-grained, se-	1574
1520	<i>2021, virtual.</i>	mantic code change patterns . In <i>Proceedings of the</i>	1575
1521	Rabee Sohail Malik, Jibesh Patra, and Michael Pradel.	<i>41st International Conference on Software Engineer-</i>	1576
1522	2019. NI2type: inferring javascript function types	<i>ing, ICSE 2019, Montreal, QC, Canada, May 25-31,</i>	1577
1523	from natural language information . In <i>Proceedings</i>	2019, pages 819–830. IEEE / ACM.	1578
1524	<i>of the 41st International Conference on Software En-</i>	Daniel Nichols, Joshua Hoke Davis, Zhaojun Xie, Ar-	1579
1525	<i>gineering, ICSE 2019, Montreal, QC, Canada, May</i>	jun Rajaram, and Abhinav Bhatel. 2024. Can large	1580
1526	<i>25-31, 2019, pages 304–315. IEEE / ACM.</i>	language models write parallel code? In <i>Proceed-</i>	1581
1527	Vishal Manjunatha Malode. 2024. <i>Benchmarking pub-</i>	ings of the 33rd International Symposium on High-	1582
1528	<i>lic large language model</i> . Ph.D. thesis, Technische	Performance Parallel and Distributed Computing,	1583
1529	Hochschule Ingolstadt.	<i>HPDC 2024, Pisa, Italy, June 3-7, 2024, pages 281–</i>	1584
1530	Timothy R McIntosh, Teo Susnjak, Nalin Arachchilage,	294. ACM.	1585
1531	Tong Liu, Dan Xu, Paul Watters, and Malka N Halga-	Pengyu Nie, Rahul Banerjee, Junyi Jessy Li, Raymond J.	1586
1532	muge. 2025. Inadequacies of large language model	Mooney, and Milos Gligoric. 2023. Learning deep	1587
1533	benchmarks in the era of generative artificial intelli-	semantics for test completion . In <i>45th IEEE/ACM</i>	1588
1534	gence. <i>IEEE Transactions on Artificial Intelligence</i> .	<i>International Conference on Software Engineering,</i>	1589
1535	Courtney Miller, Sophie Cohen, Daniel Klug, Bogdan	<i>ICSE 2023, Melbourne, Australia, May 14-20, 2023,</i>	1590
1536	Vasilescu, and Christian KaUstner. 2022. "did you	pages 2111–2123. IEEE.	1591
1537	miss my comment or what?": understanding toxicity	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan	1592
1538	in open source discussions . In <i>Proceedings of the</i>	Wang, Yingbo Zhou, Silvio Savarese, and Caiming	1593
1539	<i>44th International Conference on Software Engineer-</i>	Xiong. 2022. Codegen: An open large language	1594
1540	<i>ing, ICSE ’22, page 710–722, New York, NY, USA.</i>	model for code with multi-turn program synthesis.	1595
1541	Association for Computing Machinery.	<i>arXiv preprint arXiv:2203.13474.</i>	1596
1542	Amir M. Mir, Evaldas Latoskinas, Sebastian Proksch,	Georgios Nikitopoulos, Konstantina Dritsa, Panos	1597
1543	and Georgios Gousios. 2022. Type4py: Practical	Louridas, and Dimitris Mitropoulos. 2021. Crossvul:	1598
		a cross-language vulnerability dataset with commit	1599

1712	commit message generation . In <i>IEEE International Conference on Software Analysis, Evolution and Reengineering, SANER 2024, Rovaniemi, Finland, March 12-15, 2024</i> , pages 728–739. IEEE.	1768
1713		1769
1714		1770
1715		1771
1716	Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2024. An empirical evaluation of using large language models for automated unit test generation . <i>IEEE Transactions on Software Engineering</i> , 50(1):85–105.	1772
1717		1773
1718		1774
1719		1775
1720		1776
1721	Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, Gongye Liu, Xiaomei Nie, Deng Cai, and Yujiu Yang. 2024a. Chartmimic: Evaluating lmm’s cross-modal reasoning capability via chart-to-code generation . <i>CoRR</i> , abs/2406.09961.	1777
1722		1778
1723		1779
1724		1780
1725		1781
1726		1782
1727	Quan Shi, Michael Tang, Karthik Narasimhan, and Shunyu Yao. 2024b. Can language models solve olympiad programming? <i>CoRR</i> , abs/2404.10952.	1783
1728		1784
1729		1785
1730		1786
1731	Tianze Shi, Chen Zhao, Jordan L. Boyd-Graber, Hal Daumé III, and Lillian Lee. 2020. On the potential of lexico-logical alignments for semantic parsing to SQL queries . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020</i> , volume EMNLP 2020 of <i>Findings of ACL</i> , pages 1849–1864. Association for Computational Linguistics.	1787
1732		1788
1733		1789
1734		1790
1735		1791
1736		1792
1737		1793
1738	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1794
1739		1795
1740		1796
1741		1797
1742		1798
1743		1799
1744		1800
1745	Disha Shrivastava, Denis Kocetkov, Harm de Vries, Dzmitry Bahdanau, and Torsten Scholak. 2023a. Re-pofusion: Training code models to understand your repository . <i>CoRR</i> , abs/2306.10998.	1801
1746		1802
1747		1803
1748		1804
1749	Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. 2023b. Repository-level prompt generation for large language models of code . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 31693–31715. PMLR.	1805
1750		1806
1751		1807
1752		1808
1753		1809
1754		1810
1755		1811
1756	Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob R. Gardner, Yiming Yang, Milad Hashemi, Graham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. 2024. Learning performance-improving code edits . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1812
1757		1813
1758		1814
1759		1815
1760		1816
1761		1817
1762		1818
1763		1819
1764	Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2code: How far are we from automating front-end engineering? <i>CoRR</i> , abs/2403.03163.	1820
1765		1821
1766		1822
1767		1823
		1824
	Manav Singhal, Tushar Aggarwal, Abhijeet Awasthi, Nagarajan Natarajan, and Aditya Kanade. 2024. Nofuneval: Funny how code lms falter on requirements beyond functional correctness . <i>CoRR</i> , abs/2401.15963.	
	Patrick Suppes, Joseph L Zinnes, et al. 1962. <i>Basic measurement theory</i> .	
	Jeffrey Svajlenko, Judith F. Islam, Iman Keivanloo, Chanchal Kumar Roy, and Mohammad Mamun Mia. 2014. Towards a big data curated benchmark of inter-project code clones . In <i>30th IEEE International Conference on Software Maintenance and Evolution, Victoria, BC, Canada, September 29 - October 3, 2014</i> , pages 476–480. IEEE Computer Society.	
	Xiangru Tang, Bill Qian, Rick Gao, Jiakang Chen, Xinyun Chen, and Mark B Gerstein. 2024. Biocoder: a benchmark for bioinformatics code generation with large language models. <i>Bioinformatics</i> , 40(Supplement_1):i266–i276.	
	Wei Tao, Yanlin Wang, Ensheng Shi, Lun Du, Shi Han, Hongyu Zhang, Dongmei Zhang, and Wenqiang Zhang. 2022. A large-scale empirical study of commit message generation: models, datasets and evaluation . <i>Empir. Softw. Eng.</i> , 27(7):198.	
	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Yinxu Pan, Yesai Wu, Haotian Hui, Weichuan Liu, Zhiyuan Liu, and Maosong Sun. 2024. Debugbench: Evaluating debugging capability of large language models . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 4173–4198. Association for Computational Linguistics.	
	Michele Tufano, Shao Kun Deng, Neel Sundaresan, and Alexey Svyatkovskiy. 2022. METHODS2TEST: A dataset of focal methods mapped to test cases . In <i>19th IEEE/ACM International Conference on Mining Software Repositories, MSR 2022, Pittsburgh, PA, USA, May 23-24, 2022</i> , pages 299–303. ACM.	
	Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin White, and Denys Poshyvanyk. 2019. An empirical study on learning bug-fixing patches in the wild via neural machine translation . <i>ACM Trans. Softw. Eng. Methodol.</i> , 28(4):19:1–19:29.	
	Prashanth Vijayaraghavan, Luyao Shi, Stefano Ambrogio, Charles Mackin, Apoorva Nitsure, David Beymer, and Ehsan Degan. 2024. Vhdl-eval: A framework for evaluating large language models in VHDL code generation . <i>CoRR</i> , abs/2406.04379.	
	Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S. Yu. 2018. Improving automatic source code summarization via deep reinforcement learning . In <i>Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018, Montpellier, France, September 3-7, 2018</i> , pages 397–407. ACM.	

- Junjie Wang, Yuchao Huang, Chunyang Chen, Zhe Liu, Song Wang, and Qing Wang. 2024a. Software testing with large language models: Survey, landscape, and vision. *IEEE Transactions on Software Engineering*.
- Ping Wang, Tian Shi, and Chandan K. Reddy. 2020. Text-to-sql generation for question answering on electronic medical records. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, pages 350–361. ACM / IW3C2.
- Shuai Wang, Liang Ding, Li Shen, Yong Luo, Bo Du, and Dacheng Tao. 2024b. [OOP: object-oriented programming evaluation benchmark for large language models](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 13619–13639. Association for Computational Linguistics.
- Wenhan Wang, Chenyuan Yang, Zhijie Wang, Yuheng Huang, Zhaoyang Chu, Da Song, Lingming Zhang, An Ran Chen, and Lei Ma. 2024c. [TESTEVAL: benchmarking large language models for test case generation](#). *CoRR*, abs/2406.04531.
- Yibo Wang, Ying Wang, Tingwei Zhang, Yue Yu, Shing-Chi Cheung, Hai Yu, and Zhiliang Zhu. 2023a. [Can machine learning pipelines be better configured?](#) In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023*, page 463–475, New York, NY, USA. Association for Computing Machinery.
- Yu Emma Wang, Gu-Yeon Wei, and David Brooks. 2019. [Benchmarking tpu, gpu, and cpu platforms for deep learning](#). *Preprint*, arXiv:1907.10701.
- Zhiruo Wang, Grace Cuenca, Shuyan Zhou, Frank F. Xu, and Graham Neubig. 2023b. [Mconala: A benchmark for code generation from multiple natural languages](#). In *Findings of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 265–273. Association for Computational Linguistics.
- Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. 2022. Execution-based evaluation for open-domain code generation. *arXiv preprint arXiv:2212.10481*.
- Zora Zhiruo Wang, Akari Asai, Xinyan Velocity Yu, Frank F. Xu, Yiqing Xie, Graham Neubig, and Daniel Fried. 2024d. [Coderag-bench: Can retrieval augment code generation?](#) *CoRR*, abs/2406.14497.
- Cody Watson, Michele Tufano, Kevin Moran, Gabriele Bavota, and Denys Poshyvanyk. 2020. [On learning meaningful assert statements for unit test cases](#). In *ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020*, pages 1398–1409. ACM.
- Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2024. Jailbroken: How does llm safety training fail? *Advances in Neural Information Processing Systems*, 36.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*.
- Jiayi Wei, Greg Durrett, and Isil Dillig. 2023. [Typet5: Seq2seq type inference using static analysis](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo. 2024a. [Plot2code: A comprehensive benchmark for evaluating multi-modal large language models in code generation from scientific plots](#). *CoRR*, abs/2405.07990.
- Minghao Wu, Weixuan Wang, Sinuo Liu, Huifeng Yin, Xintong Wang, Yu Zhao, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2025. [The bitter lesson learned from 2,000+ multilingual benchmarks](#). *Preprint*, arXiv:2504.15521.
- Tongtong Wu, Weigang Wu, Xingyu Wang, Kang Xu, Suyu Ma, Bo Jiang, Ping Yang, Zhenchang Xing, Yuan-Fang Li, and Gholamreza Haffari. 2024b. [Versi-code: Towards version-controllable code generation](#). *CoRR*, abs/2406.07411.
- Chunqiu Steven Xia, Yinlin Deng, and Lingming Zhang. 2024a. [Top leaderboard ranking = top coding proficiency, always? evoeval: Evolving coding benchmarks via LLM](#). *CoRR*, abs/2403.19114.
- Yinghui Xia, Yuyan Chen, Tianyu Shi, Jun Wang, and Jinsong Yang. 2024b. [Aicodereval: Improving AI domain code generation of large language models](#). *CoRR*, abs/2406.04712.
- Jie Xiao, Qianyi Huang, Xu Chen, and Chen Tian. 2024. [Large language model performance benchmarking on mobile platforms: A thorough evaluation](#). *Preprint*, arXiv:2410.03613.
- Ruiyang Xu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xi-anpei Han, Ben He, Shing-Chi Cheung, and Le Sun. 2024. [Cruxeval-x: A benchmark for multilingual code reasoning, understanding and execution](#). *CoRR*, abs/2408.13001.
- Ankit Yadav, Himanshu Beniwal, and Mayank Singh. 2024a. [Pythonsaga: Redefining the benchmark to evaluate code generating llms](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 17113–17126.
- Ankit Yadav, Himanshu Beniwal, and Mayank Singh. 2024b. [Pythonsaga: Redefining the benchmark to evaluate code generating llms](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 17113–17126. Association for Computational Linguistics.

1937	Shuhan Yan, Hang Yu, Yuting Chen, Beijun Shen, and Lingxiao Jiang. 2020. Are the code snippets what we are searching for? A benchmark and an empirical study on code search with natural-language queries . In <i>27th IEEE International Conference on Software Analysis, Evolution and Reengineering, SANER 2020, London, ON, Canada, February 18-21, 2020</i> , pages 344–354. IEEE.	Hao Yu, Bo Shen, Dezhi Ran, Jiaxin Zhang, Qi Zhang, Yuchi Ma, Guangtai Liang, Ying Li, Tao Xie, and Qianxiang Wang. 2023. Codereval: A benchmark of pragmatic code generation with generative pre-trained models . <i>arXiv preprint arXiv:2302.00288</i> .	1995
1938			1996
1939			1997
1940			1998
1941			1999
1942			
1943			
1944			
1945	Weixiang Yan, Yuchen Tian, Yunzhe Li, Qian Chen, and Wen Wang. 2023. Codetransocean: A comprehensive multilingual benchmark for code translation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 5067–5089. Association for Computational Linguistics.	Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander R. Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Walter S. Lasecki, and Dragomir R. Radev. 2019a. Cosql: A conversational text-to-sql challenge towards cross-domain natural language interfaces to databases . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 1962–1979. Association for Computational Linguistics.	2000
1946			2001
1947			2002
1948			2003
1949			2004
1950			2005
1951			2006
			2007
			2008
			2009
1952	Hongyu Yang, Liyang He, Min Hou, Shuanghong Shen, Rui Li, Jiahui Hou, Jianhui Ma, and Junda Zhao. 2024a. Aligning llms through multi-perspective user preference ranking-based feedback for programming question answering . <i>CoRR</i> , abs/2406.00037.		2010
1953			2011
1954			2012
1955			2013
1956			2014
1957	Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong, Xu Han, Yukun Yan, Zhenghao Liu, Zhixing Tan, Pengyuan Liu, Dong Yu, Zhiyuan Liu, Xiaodong Shi, and Maosong Sun. 2024b. Matplotagent: Method and evaluation for llm-based agentic scientific data visualization . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 11789–11804. Association for Computational Linguistics.	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 3911–3921. Association for Computational Linguistics.	2015
1958			2016
1959			2017
1960			2018
1961			2019
1962			2020
1963			2021
1964			2022
1965			2023
1966			2024
1967	Ziyu Yao, Daniel S. Weld, Wei-Peng Chen, and Huan Sun. 2018. Staqa: A systematically mined question-code dataset from stack overflow . In <i>Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018</i> , pages 1693–1703. ACM.	Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir R. Radev. 2019b. Sparc: Cross-domain semantic parsing in context . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 4511–4523. Association for Computational Linguistics.	2025
1968			2026
1969			2027
1970			2028
1971			2029
1972			2030
			2031
1973	Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models . <i>Preprint</i> , arXiv:2303.10420.		2032
1974			2033
1975			2034
1976			2035
1977			
1978			
1979	Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. 2018. Learning to mine aligned code and natural language pairs from stack overflow . In <i>International Conference on Mining Software Repositories, MSR</i> , pages 476–486. ACM.	Xiao Yu, Lei Liu, Xing Hu, Jacky Wai Keung, Jin Liu, and Xin Xia. 2024. Fight fire with fire: How much can we trust chatgpt on source code-related tasks? <i>arXiv preprint arXiv:2405.12641</i> .	2036
1980			2037
1981			2038
1982			2039
1983			
1984			
1985	Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Oleksandr Polozov, and Charles Sutton. 2023. Natural language to code generation in interactive data science notebooks . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 126–173. Association for Computational Linguistics.	Xiaoqing Yu, Tianlong Chen, Zhengjie Yu, Huiyu Li, Yang Yang, Xiaoqian Jiang, and Anxiao Jiang. 2020. Dataset and enhanced model for eligibility criteria-to-sql semantic parsing . In <i>Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020</i> , pages 5829–5837. European Language Resources Association.	2040
1986			2041
1987			2042
1988			2043
1989			2044
1990			2045
1991			2046
1992			2047
1993			
1994			
		Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. 2024a. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher . <i>Preprint</i> , arXiv:2308.06463.	2048
			2049
			2050
			2051

2052	Zhiqiang Yuan, Mingwei Liu, Shiji Ding, Kaixin Wang,	13643–13658. Association for Computational Lin-	2109
2053	Yixuan Chen, Xin Peng, and Yiling Lou. 2024b. Eval-	guistics.	2110
2054	uating and improving chatgpt for unit test generation.		
2055	<i>Proceedings of the ACM on Software Engineering</i> ,	Quanjun Zhang, Tongke Zhang, Juan Zhai, Chunrong	2111
2056	1(FSE):1703–1726.	Fang, Bowen Yu, Weisong Sun, and Zhenyu Chen.	2112
		2023a. A critical review of large language model on	2113
2057	Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding,	software engineering: An example from chatgpt and	2114
2058	Kaixin Wang, Yixuan Chen, and Xin Peng. 2023a.	automated program repair . <i>CoRR</i> , abs/2310.08879.	2115
2059	No more manual tests? evaluating and improving		
2060	chatgpt for unit test generation. <i>arXiv preprint</i>	Shudan Zhang, Hanlin Zhao, Xiao Liu, Qinkai Zheng,	2116
2061	<i>arXiv:2305.04207</i> .	Zehan Qi, Xiaotao Gu, Xiaohan Zhang, Yuxiao Dong,	2117
		and Jie Tang. 2024b. Naturalcodebench: Examining	2118
2062	Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji	coding performance mismatch on humaneval and	2119
2063	Ding, Kaixin Wang, Yixuan Chen, and Xin Peng.	natural user prompts . <i>CoRR</i> , abs/2405.04520.	2120
2064	2023b. No more manual tests? evaluating and		
2065	improving chatgpt for unit test generation . <i>CoRR</i> ,	Yusen Zhang, Jun Wang, Zhiguo Wang, and Rui Zhang.	2121
2066	abs/2305.04207.	2023b. Xsemplr: Cross-lingual semantic parsing in	2122
		multiple natural languages and meaning representa-	2123
2067	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng,	tions . In <i>Proceedings of the 61st Annual Meeting of</i>	2124
2068	Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu	<i>the Association for Computational Linguistics (Vol-</i>	2125
2069	Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao	<i>ume 1: Long Papers)</i> , <i>ACL 2023, Toronto, Canada,</i>	2126
2070	Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan	<i>July 9-14, 2023</i> , pages 15918–15947. Association for	2127
2071	Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang,	Computational Linguistics.	2128
2072	Huan Sun, Yu Su, and Wenhao Chen. 2024. Mmmu:		
2073	A massive multi-discipline multimodal understand-	Ziyin Zhang, Lizhen Xu, Zhaokun Jiang, Hongkun	2129
2074	ing and reasoning benchmark for expert agi . <i>Preprint</i> ,	Hao, and Rui Wang. 2024c. Multiple-choice ques-	2130
2075	arXiv:2311.16502.	tions are efficient and robust LLM evaluators . <i>CoRR</i> ,	2131
		abs/2405.11966.	2132
2076	Sukmin Yun, Haokun Lin, Rusiru Thushara, Moham-		
2077	mad Qazim Bhat, Yongxin Wang, Zutao Jiang,	Dewu Zheng, Yanlin Wang, Ensheng Shi, Ruikai	2133
2078	Mingkai Deng, Jinhong Wang, Tianhua Tao, Junbo	Zhang, Yuchi Ma, Hongyu Zhang, and Zibin Zheng.	2134
2079	Li, Haonan Li, Preslav Nakov, Timothy Baldwin,	2024. Towards more realistic evaluation of llm-based	2135
2080	Zhengzhong Liu, Eric P. Xing, Xiaodan Liang, and	code generation: an experimental study and beyond .	2136
2081	Zhiqiang Shen. 2024. Web2code: A large-scale	<i>CoRR</i> , abs/2406.06918.	2137
2082	webpage-to-code dataset and evaluation framework		
2083	for multimodal llms . <i>CoRR</i> , abs/2406.20098.	Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan	2138
		Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang,	2139
2084	Mohamed Zairi. 2010. <i>Benchmarking for best practice</i> .	Yang Li, Teng Su, Zhilin Yang, and Jie Tang. 2023a.	2140
2085	Routledge.	Codegeex: A pre-trained model for code generation	2141
		with multilingual benchmarking on humaneval-x . In	2142
2086	Daoguang Zan, Bei Chen, Zeqi Lin, Bei Guan, Yongji	<i>Proceedings of the 29th ACM SIGKDD Conference</i>	2143
2087	Wang, and Jian-Guang Lou. 2022a. When lan-	<i>on Knowledge Discovery and Data Mining, KDD '23,</i>	2144
2088	guage model meets private library. <i>arXiv preprint</i>	page 5673–5684, New York, NY, USA. Association	2145
2089	<i>arXiv:2210.17236</i> .	for Computing Machinery.	2146
2090	Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin,	Yunhui Zheng, Saurabh Pujar, Burn L. Lewis, Luca	2147
2091	Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen,	Buratti, Edward A. Epstein, Bo Yang, Jim Laredo,	2148
2092	and Jian-Guang Lou. 2022b. Cert: continual pre-	Alessandro Morari, and Zhong Su. 2021. D2A: A	2149
2093	training on sketches for library-oriented code genera-	dataset built for ai-based vulnerability detection meth-	2150
2094	tion. <i>arXiv preprint arXiv:2206.06888</i> .	ods using differential analysis . In <i>43rd IEEE/ACM</i>	2151
		<i>International Conference on Software Engineering:</i>	2152
2095	Zhengran Zeng, Yidong Wang, Rui Xie, Wei Ye, and	<i>Software Engineering in Practice, ICSE (SEIP) 2021,</i>	2153
2096	Shikun Zhang. 2024. Coderujb: An executable and	<i>Madrid, Spain, May 25-28, 2021</i> , pages 111–120.	2154
2097	unified java benchmark for practical programming	IEEE.	2155
2098	scenarios . In <i>Proceedings of the 33rd ACM SIGSOFT</i>		
2099	<i>International Symposium on Software Testing and</i>	Zibin Zheng, Kaiwen Ning, Yanlin Wang, Jingwen	2156
2100	<i>Analysis, ISSTA 2024, Vienna, Austria, September</i>	Zhang, Dewu Zheng, Mingxi Ye, and Jiachi Chen.	2157
2101	<i>16-20, 2024</i> , pages 124–136. ACM.	2023b. A survey of large language models for code:	2158
		Evolution, benchmarking, and future trends. <i>arXiv</i>	2159
2102	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin.	<i>preprint arXiv:2311.10372</i> .	2160
2103	2024a. Codeagent: Enhancing code generation with		
2104	tool-integrated agent systems for real-world repo-	Victor Zhong, Caiming Xiong, and Richard Socher.	2161
2105	level coding challenges . In <i>Proceedings of the 62nd</i>	2017. Seq2sql: Generating structured queries	2162
2106	<i>Annual Meeting of the Association for Computa-</i>	from natural language using reinforcement learning .	2163
2107	<i>tional Linguistics (Volume 1: Long Papers)</i> , <i>ACL</i>	<i>CoRR</i> , abs/1709.00103.	2164
2108	<i>2024, Bangkok, Thailand, August 11-16, 2024</i> , pages		

- Yaqin Zhou, Shangqing Liu, Jing Kai Siow, Xiaoning Du, and Yang Liu. 2019. [Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 10197–10207.
- Ming Zhu, Aneesh Jain, Karthik Suresh, Roshan Ravindran, Sindhu Tipirneni, and Chandan K. Reddy. 2022. [Xlcost: A benchmark dataset for cross-lingual code intelligence](#). *CoRR*, abs/2206.08474.
- Qiming Zhu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Shing-Chi Cheung. 2024. [Domaineval: An auto-constructed benchmark for multi-domain code generation](#). *CoRR*, abs/2408.13204.
- Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kadour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro von Werra. 2024. [Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions](#). *CoRR*, abs/2406.15877.
- Deqing Zou, Sujuan Wang, Shouhuai Xu, Zhen Li, and Hai Jin. 2020. [\$\mu\$ vuldeepecker: A deep learning-based system for multiclass vulnerability detection](#). *CoRR*, abs/2001.02334.

A Detailed Guidance in HOW2BENCH

We present the detailed guidance in each phase in this section.

A.1 Guideline for Benchmark Design

	Phase 1. Benchmark Design	Priority	☑
1	Consider whether the benchmark can fill the gap in related research.	★★★	<input type="checkbox"/>
2	Consider what is the expected scope of the benchmark set (e.g., what natural languages, programming languages, task granularity).	★★★	<input type="checkbox"/>
3	Consider the expected application scenario of this benchmark (e.g., programming assistant, automated tester).	★★★	<input type="checkbox"/>
4	Consider the LLMs' capabilities (e.g., understanding, reasoning, calculation) and domain knowledge (e.g., OOP, memory management, fault localization, process scheduling) that the benchmark hopes to evaluate.	★★★	<input type="checkbox"/>

Figure 3: Guideline for Benchmark Design

A.2 Guideline for Benchmark Evaluation

	Phase 3. Benchmark Evaluation	Priority	☑
24	Consider whether sufficient LLMs are evaluated.	★	<input type="checkbox"/>
25	Consider whether representative LLMs (e.g., covering latest/classical LLM families, small/large LLMs, and open/closed-source LLMs) are evaluated.	★	<input type="checkbox"/>
26	Consider whether the prompt is of high quality (e.g., the instruction and intent are clear).	★★★	<input type="checkbox"/>
27	The prompts have been validated by humans or LLMs (e.g., evaluated or discussed by participants or preliminarily tried out on several LLMs).	★	<input type="checkbox"/>
28	Try different paraphrases of the prompt.	★	<input type="checkbox"/>
29	Try different prompting strategies to observe the impact on the evaluation results (e.g., in-context learning, chain-of-thought).	★★	<input type="checkbox"/>
30	Pay attention to the hardware environment (such as GPU card, storage size, etc.) of the experiment.	★★★	<input type="checkbox"/>
31	Pay attention to the operating system and software environment (e.g. operating system, version, etc.) used for the experiment.	★★★	<input type="checkbox"/>
32	Pay attention to the off-the-shelf platforms, frameworks, or libraries for LLM evaluation (e.g., fast chat, vllm, huggingface) that are used.	★★	<input type="checkbox"/>
33	Repeat the experiment multiple times to reduce the impact of randomness on the evaluation.	★	<input type="checkbox"/>
34	Consider various randomization strategies (e.g., trying various temperature parameters) to reduce the impact of parameter configuration on the evaluation.	★★	<input type="checkbox"/>
35	Record the experimental process in detail (e.g., parameter settings, running time, input/output pairs, etc.).	★★★	<input type="checkbox"/>

Figure 4: Guideline for Benchmark Evaluation

A.3 Guideline for Benchmark Construction

	Phase 2. Benchmark Construction (1/2)	Priority	☑
5	Consider whether the data source of the benchmark is traceable.	★	<input type="checkbox"/>
6	Consider whether the data source of the benchmark is of high quality (e.g., stars, downloads, last update times, number of forks).	★★	<input type="checkbox"/>
7	Consider whether the benchmark's data source is representative (e.g., choose an open-source community or code hosting platform that matches the task, capability, and scope under study).	★	<input type="checkbox"/>
8	Consider data contamination issues during the benchmark collection (e.g., considering the upload time of the source code or checking whether the data source is included in the training data of LLMs).	★	<input type="checkbox"/>
9	If data sampling is needed, consider whether the choice of sample size is scientific (e.g., considering the confidence level/margin of error/sampling proportion, etc.).	★	<input type="checkbox"/>
10	If data sampling is needed, consider whether the sampling process is rigorous (e.g., random sampling, stratified sampling, etc.).	★	<input type="checkbox"/>
11	Ensure each data point in the benchmark falls into the targeted scope (e.g., checking each data point's evaluated capabilities or domain knowledge).	★	<input type="checkbox"/>
12	Consider whether the data in the benchmark can cover the studied capabilities/domain knowledge/application scenarios.	★★★	<input type="checkbox"/>
13	Consider whether there is a standard answer for each sample in the benchmark (such as reference code, etc.).	★★	<input type="checkbox"/>
14	For code, consider whether the code is compilable/executable.	★	<input type="checkbox"/>
15	Consider the possibility of noise in the data and perform denoise.	★★	<input type="checkbox"/>
16	Consider the possibility of duplication in the data and deduplicate them.	★★	<input type="checkbox"/>
17	Clean the sensitive information (such as data desensitization and anonymization) unless the benchmark is deliberately designed so.	★★	<input type="checkbox"/>
18	Manually review some or all of the data in the benchmark to ensure its quality.	★★	<input type="checkbox"/>
19	Use LLMs to review some or all of the data in the benchmark to ensure its quality.	★	<input type="checkbox"/>
20	Design appropriate output validation methods for the benchmark (e.g., using exact matching or designing test cases).	★★★	<input type="checkbox"/>
21	Design appropriate evaluation metrics for the evaluation set (e.g., precision, accuracy, pass@K, recall).	★★★	<input type="checkbox"/>
22	Consider the adequacy of the evaluation metrics (e.g., is the code coverage high enough).	★★★	<input type="checkbox"/>
23	Consider if there are any other evaluation perspectives (e.g., readability, efficiency, safety, security).	★	<input type="checkbox"/>

Figure 5: Guideline for Benchmark Construction

A.4 Guideline for Benchmark Analysis

	Phase 4. Benchmark Analysis	Priority	☑
36	Observe the difficulty of the benchmark, checking if the benchmark is too hard or too easy for LLMs (i.e., most LLMs score too high/low).	★★	<input type="checkbox"/>
37	Consider whether the benchmark can distinguish the pros and cons of different LLMs.	★	<input type="checkbox"/>
38	If the experiment is repeated several times, consider the stability of the benchmark (i.e., whether the experimental results vary too much in the repeated experiments).	★	<input type="checkbox"/>
39	Analyze the correlation between the data and their score. For example, if there is a correlation between the data (such as similar difficulty and knowledge required), then the scores should also be correlated.	★★	<input type="checkbox"/>
40	Compare the performance of LLMs on this benchmark with their performance on other related benchmarks.	★	<input type="checkbox"/>
41	Consider presenting the experiment results in an appropriate way (e.g., table, line graph, pie chart, etc.).	★★★	<input type="checkbox"/>
42	Consider presenting the experiment results clearly (e.g., distinguishable colors/labels/shapes, etc.).	★★★	<input type="checkbox"/>
43	Explain the experiment results.	★★★	<input type="checkbox"/>
44	Observe correlations via multiple perspectives from the experimental results (e.g., performance is correlated with model size or amount of context).	★	<input type="checkbox"/>

Figure 6: Guideline for Evaluation Analysis

A.5 Guideline for Benchmark Release

	Phase 5. Benchmark Release	Priority	☑
45	The analysis of the evaluation results will be inspiring (e.g., shed light on future direction, make actionable advice, etc.).	★	<input type="checkbox"/>
46	Set the appropriate license for the benchmark.	★★★	<input type="checkbox"/>
47	Review the released benchmark or other artifacts to ensure they do NOT contain sensitive information (e.g., API keys, usernames, passwords, etc.).	★★★	<input type="checkbox"/>
48	review the released benchmark or other artifacts to ensure they do NOT contain toxicity information (e.g., abusive comments/identifiers).	★★★	<input type="checkbox"/>
49	Make sure the benchmark is open-accessible.	★★★	<input type="checkbox"/>
50	Make sure the test cases or reference data are open and accessible.	★★★	<input type="checkbox"/>
51	Provide prompts used in the experiment to ensure the experiments are reproducible.	★★★	<input type="checkbox"/>
52	Disclose the experimental environment (e.g., hardware, operating system, software version, framework platform) to ensure the reproducibility of the experiment.	★★★	<input type="checkbox"/>
53	Make the detailed experimental results public for verification.	★★★	<input type="checkbox"/>
54	Ensure the quality of the user manual such as README (e.g., it contains necessary benchmark introduction, executable scripts, etc.).	★★	<input type="checkbox"/>
55	Provide convenient evaluation interfaces for the released benchmark (e.g., providing a command line interface, docker, etc.).	★★	<input type="checkbox"/>

Figure 7: Guideline for Benchmark Release

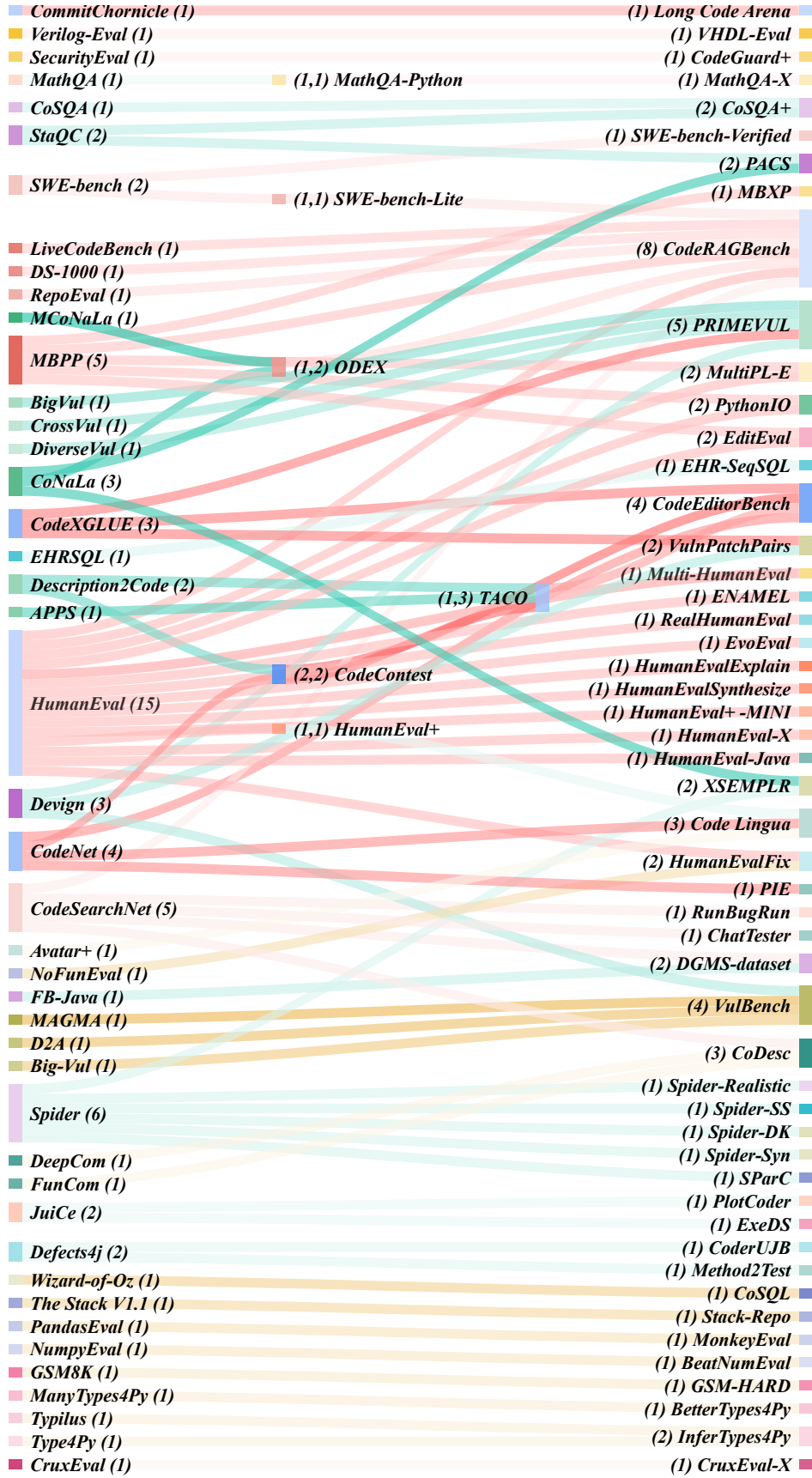


Figure 8: Relationships between Benchmarks

B Statistics of studied benchmarks

In this section, we conducted a comprehensive and detailed statistical analysis of the 274 benchmarks collected.

B.1 Profile of Studied Benchmarks

We first show the trend in the development of benchmarks from 2014 to 2024. As shown in Figure 9, the data shows a modest beginning, with only a handful of benchmarks created annually until 2017. From 2018 onwards, there is a noticeable uptrend in benchmark creation, culminating in a significant jump to 149 benchmarks in 2024. This sharp increase indicates a recent heightened interest and demand for comprehensive code-related benchmarks for LLMs, reflecting the evolving complexities and expanding requirements of automated software engineering.

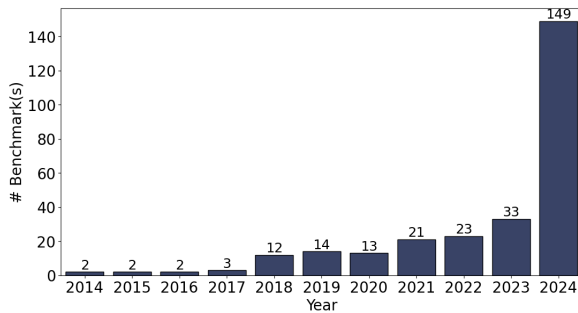


Figure 9: Benchmark Distribution over Years

Hierarchy of Benchmarks. Figure 8 visualizes the inheritance relationships among benchmarks, indicating that the benchmarks on the *left serve as sources* for those on the right. It highlights that **18% (50 out of 274) of benchmarks act as data sources**, continuously benefiting the construction of subsequent benchmarks.

Figure 8 reveals that *HumanEval* (Chen et al., 2021a), as the **most significant source** benchmark, benefits at least **15 downstream benchmarks**, followed by MBPP (Austin et al., 2021) and CodeSearchNet (Husain et al., 2019). From the right side of the figure, some benchmarks, like VulBench (Gao et al., 2023b), incorporate methodologies or data from 4 previous benchmarks, and codeRagBench (Wang et al., 2024d) integrates elements from 8 prior benchmarks.

This hierarchical structure among benchmarks also alerts us that the **data quality of a benchmark not only affects its own credibility but can continue to impact others** if it serves as a source. This

underscores the importance of adhering to stringent guidelines during benchmark development and highlights the crucial role of **establishing standards** to ensure the integrity and utility of benchmark data across research and development efforts.

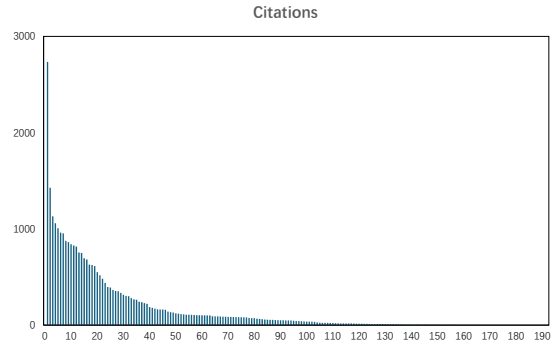


Figure 10: Citation Distribution of Benchmarks

Annual Trend. Regarding the *coding tasks*, Figure 11 illustrates the distribution of various coding tasks across benchmarks. It is clear that the task of **Code Generation** is most prevalent, with 99 benchmarks focusing on this area, according to 36% (99/274) of studied benchmarks, indicating a significant interest in generating code automatically. Program Repair and Defect Detection are well-represented, with 27 and 25 benchmarks, respectively, reflecting the importance of correcting code and detecting defects.

Citation distribution. We also visualized the citations of 274 code-related benchmarks. The citation statistics were collected on September 1st, 2024. From Figure 10, we can see a clear long-tail trend of the citations, from the highest 2735 (HumanEval (Chen et al., 2021a)) to the lowest 0.

Coding Task. Tasks like Code Summarization and Text2SQL are similarly significant, each with 25 and 22 benchmarks. These tasks focus on making code more understandable and converting natural language queries into SQL queries. Other tasks, such as Code Retrieval, Code Reasoning, and Code Translation, are represented with 18, 17, and 16 benchmarks, respectively. Lesser-represented benchmarks are Test Generation, Code Optimization, and Code Completion, each represented by 8 and 7 benchmarks, indicating the inadequacy of these tasks.

Programming Languages. Figure 12 shows the distribution of benchmarks across various programming languages. The overall trend indicates a strong preference for benchmarking **Python**, which

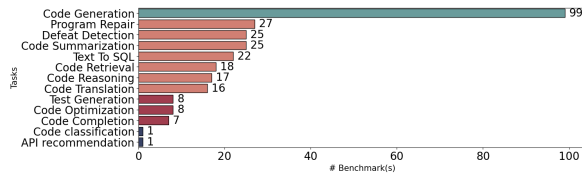


Figure 11: Benchmark Distribution over Tasks

leads with 158 benchmarks, followed by Java and C++, with 107 and 63, respectively. The graph also reveals a diverse range of languages being used. In total, 724 programming languages are studied by these 274 benchmarks. Though some programming languages, such as Kotlin, Swift, and Scala, are less frequently exercised, the benchmarks involving them are tailored to different application needs and technology environments. This distribution shows the existing benchmarks are dominated by three mainstream programming languages, leaving other programming languages less studied and benchmarked.

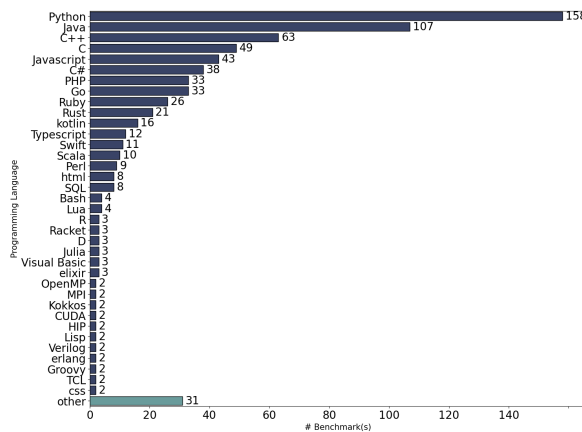


Figure 12: Benchmark Distribution over Programming Language

Natural Language. Figure 13 illustrates the distribution of benchmarks for different natural languages. The bar chart overwhelmingly shows that English is the dominant language, with 192 benchmarks highlighting its ubiquity in research and development. Other languages have significantly fewer benchmarks, with six for Chinese and only two each for Japanese, Russian, and Spanish. The category labeled “Other” includes 20 benchmarks spread across other natural languages, indicating some diversity but limited attention to non-English benchmarks. This distribution highlights the prominence of English in the global research community and also demonstrates the *uneven representation*

of natural languages in the studied benchmarks.

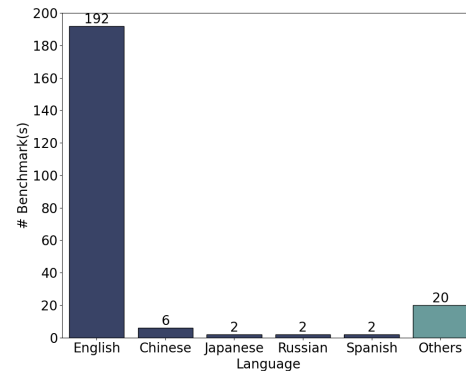


Figure 13: Benchmark Distribution over Natural Language

Modals in the benchmarks. Figure 14 presents the distribution of benchmarks according to the type of language used in their prompts. The chart shows that the majority, at 47.1%, of the benchmarks use a combination of natural language and programming Language, followed by PL only (31.0%) and NL only (21.9%).

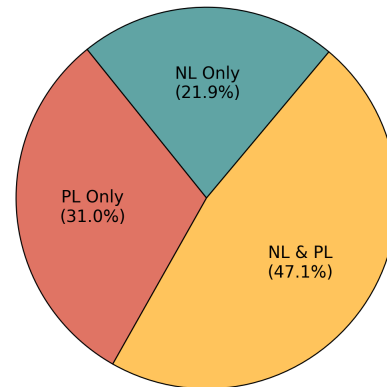


Figure 14: Benchmark Distribution over Modal in Prompt

Granularity. The code snippet in a code-related benchmark varies from statement-level (i.e., one line of code. For example, CoNaLa (Yin et al., 2018) and Math-QA (Amini et al., 2019)), function-level (i.e., a function unit of code. For example, HumanEval (Chen et al., 2021a) and MBPP (Austin et al., 2021)), class-level (i.e., a class with multiple function units of code. For example, ClassEval (Du et al., 2023b)) and project-level (i.e., a project with multiple classes or modules. For example, DevEval (Li et al., 2024a) and JavaBench (Cao et al., 2024a)). Figure 15 illustrates the granularity levels at which benchmarks are typically conducted. The chart shows that the majority of bench-

marks, comprising **71.8%, focus on the function level**. Projects constitute 15.0% of the benchmarks. Class-level granularity is the least represented at 2.6%.

The majority of benchmarks are currently at the function level (70+%), followed by the project level (15+%). This indicates that the current major demand is for **assessing individual functions within a single task**, followed by the demand for evaluating functionalities more aligned with **actual project-level code development**.

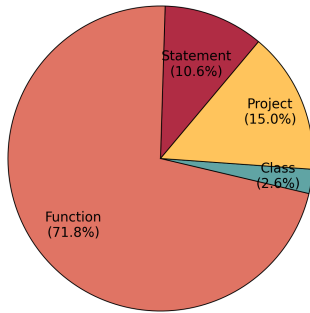


Figure 15: Benchmark Distribution over Granularity

B.2 Statistics about Benchmark Design

Design of Studied Capabilities. To understand whether benchmark developers recognize the capabilities of LLMs they aim to evaluate, we carefully analyzed 30 representative benchmarks (Appendix D) to see if they clearly specify the capabilities being assessed by their benchmarks. As shown in Figure 16, 90% of benchmarks explicitly specify the capabilities (e.g., intention understanding, problem solving, testing, debugging capabilities) to be evaluated, while 10% do not. The statistics show that the most highly cited benchmarks clearly define the assessment capabilities.

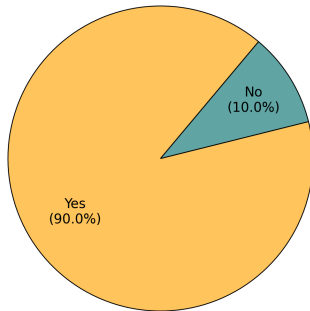


Figure 16: Benchmark Distribution Over Capabilities Consideration

Furthermore, we investigated the 30 focused

benchmarks and identified a case (Figure 17) from MBPP (Austin et al., 2021) where the case is likely to fall outside of the targeted capability of the benchmark. In particular, MBPP (Austin et al., 2021) aims to “measure the ability of these models to synthesize short Python programs from natural language descriptions” for “entry-level programmers”. As we can see from Figure 17, the prompt requires LLMs to “Write a function to calculate the dogs’ years.” Simply from this description, an entry-level programmer is unlikely to write a correct program without knowing the conversion equation from dogs’ year to dogs’ age. In other words, this case is more about assessing whether LLMs have acquired this specific knowledge rather than evaluating the most fundamental programming skills.

...
</>

Out of Targeted Capabilities

```

1 {
2   'source_file': 'Benchmark Questions
3                     Verification V2.ipynb',
4   'task_id': 264,
5   'prompt': 'Write a function to calculate
6               a dog's age in dog's years.'
7   'test_list': [ "assert dog_age(12)==61",
8                  "assert dog_age(15)==73",
9                  "assert dog_age(24)==109" ]
10 }
```

Figure 17: An Example of Out-of-capability Case from MBPP (Austin et al., 2021).

Design of Studied Application Scenarios. Similarly, to understand whether benchmark developers scoped the application scenarios of LLMs they aim to evaluate, we carefully analyzed 30 representative benchmarks (Appendix D) to see whether they explicitly specify the application scenarios their benchmarks target. As shown in Figure 18, 70% representative benchmarks have clearly specified application scenarios (e.g., programming assistant), while the rest do not. Indeed, clearly defining the application scenarios could help benchmark constructors establish precise goals for the design and development of the benchmark, ensuring accuracy in the evaluation.

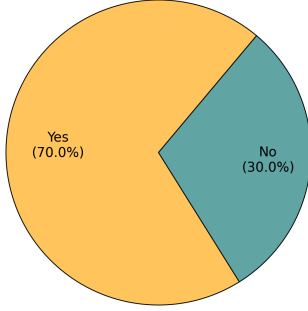


Figure 18: Benchmark Distribution Over Expected Application Scenario Consideration

B.3 Statistics about Data Preparation

B.3.1 Data Preprocessing

Data Deduplication. During benchmark preparation, data cleaning and preprocessing are necessary. However, as shown in Figure 19, only **38% benchmarks have deduplicated** the collected data. More than half of them didn’t mention this process.

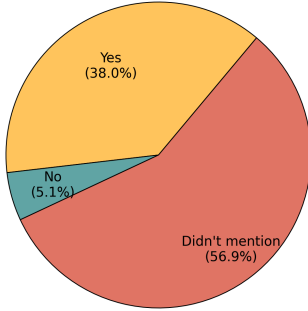


Figure 19: Benchmark Distribution over Deduplication

To investigate the situation, we went through the 30 representative benchmarks (Listed in Appendix D) and found two duplicated subjects in MBPP (Austin et al., 2021). Tasks with id 71 and 141 examined the same functionality, i.e., “Write a function to sort a list of elements.”, collected from the same source.

The significance of data preprocessing, such as **deduplication**, is frequently **overlooked** by benchmark builders, leading to data duplication even in highly cited benchmarks.

Data Quality Assurance. Ensuring data quality for the benchmark is essential. However, our statistics (Figure 21) show disappointing results. **67.9% of benchmarks do NOT take any measures for data quality assurance.** Among those benchmarks that do incorporate data quality measures, the majority rely on manual checks, which accounts for

```

1 {
2   'source_file': 'Mike's Copy of Benchmark
3     Questions Verification V2.ipynb',
4   'task_id': 71,
5   'prompt': 'Write a function to sort a list of elements.'
6 }
7
8 {
9   'source_file': 'Mike's Copy of Benchmark
10     Questions Verification V2.ipynb',
11   'task_id': 141,
12   'prompt': 'Write a function to sort a list of elements.'
13 }

```

Figure 20: A Counterexample of Rule 16 from MBPP (Austin et al., 2021).

22.6%. Other countermeasurements, such as code execution, constitute only 2.2%, while verification using LLMs accounts for 1.5%. Additional methods, such as using download counts as a basis, are also employed.

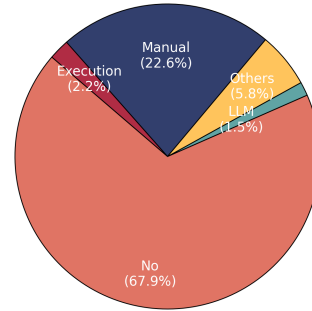


Figure 21: Benchmark Distribution over Quality Assurance Method

Additionally, we dived into the 30 representative benchmarks (Listed in Appendix D) and identified an example where the code cannot be executed successfully. As shown in Figure 22, the function `swap()` in line 7 has not been defined, so the execution of the code would fail if the code has been executed. This highlights a significant gap in ensuring the reliability and validity of benchmark data, underscoring the need for more rigorous and automated data quality assurance practices.

Data Contamination Resolution. Data contamination (Golchin and Surdeanu, 2023; Cao et al., 2024b) threat has been widely discussed. A benchmark with contaminated data may yield overclaimed results, misleading the understanding of the LLMs’ capabilities. According to our statistics (Figure 23 on benchmarks from the year

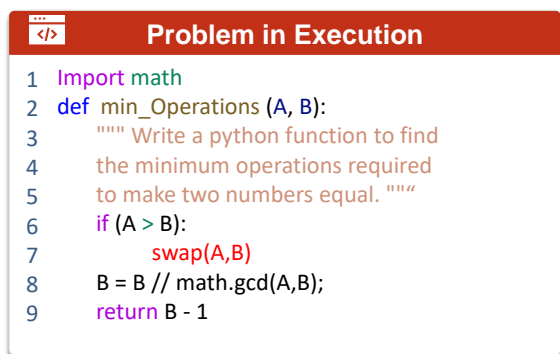


Figure 22: An Example from MBPP (Austin et al., 2021) that failed to be executed.

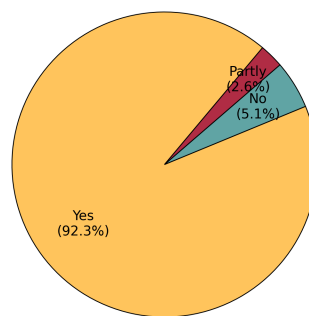


Figure 24: Benchmark Distribution over Solution

2023 to 2024 (the duration when most LLMs were launched), most (81.8 %) benchmarks were not aware of and have not taken any measures to alleviate data contamination, being vulnerable to data contamination threat.

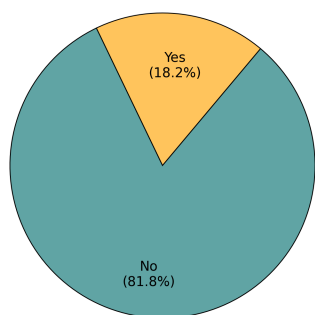


Figure 23: Benchmark Distribution over Quality Assurance on Data Contamination

B.3.2 Statistics about Data Curation

Ground truth solutions. Figure 24 shows that although the majority (92.3%) of benchmarks provide reference code as ground truth, there are 5% of benchmarks without reference code. Although it is not compulsory as long as object measurements (e.g., test cases) are provided, *a reference code is still recommended*. Indeed, if a benchmark provides reference code, its reliability tends to be better because it ensures that there are feasible solutions for the tasks involved. This guarantees that the tasks are theoretical and practically solvable, enhancing the benchmark’s usefulness and credibility in real-world applications.

Additionally, *the correctness of the ground truth solution* should also be noted. Figure 25 shows an *incorrect* code solution provided in HumanEval (Chen et al., 2021a). This should draw

benchmark constructors’ attention to the correctness of the benchmark reference code.

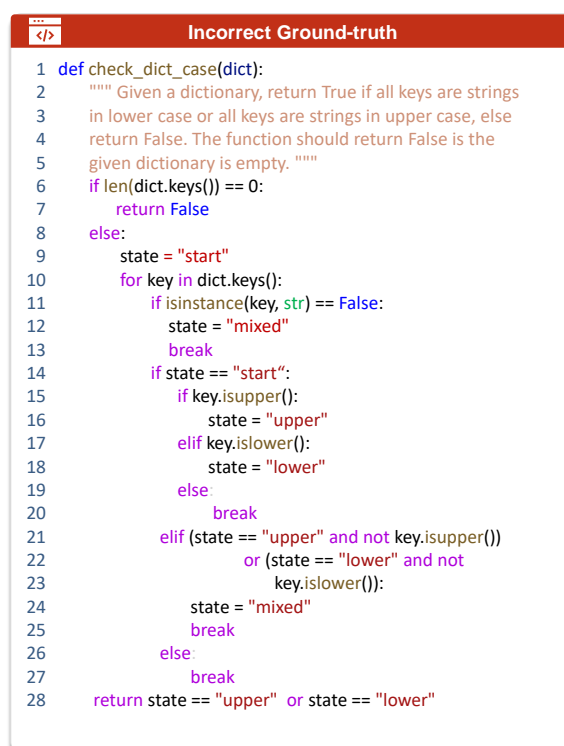


Figure 25: An Example from HumanEval (Chen et al., 2021a) which shows *an incorrect solution* provided in the benchmark.

Oracle. An oracle (Barr et al., 2014) is a way to determine whether the output is correct or not. For example, assume the output of LLMs is in the form of code, then an oracle could be running tests against the code and see whether the code can pass all the tests. Figure 26 shows the distribution of types of oracle that are used in these benchmarks. Clear that the *exact match* 41.97% (115/274) and *test case passing* (114/274) 41.6% are the most common oracle used in code-related benchmarks, followed by thresholds (i.e., similarities smaller than a specific threshold).

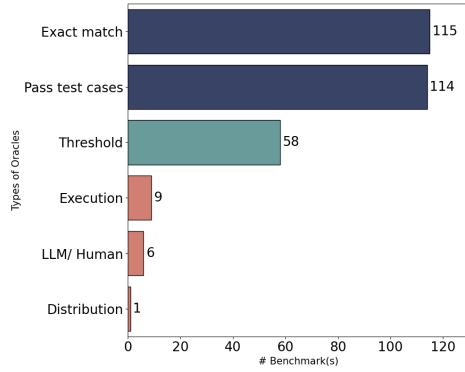


Figure 26: Benchmark Distribution over Test Oracle

Code test coverage (Ivanković et al., 2019), as a common oracle for code-related benchmarks, has been widely adopted to determine the output correctness. It should be considered if a benchmark uses test case passing as a criterion for the correctness of the generated code. Otherwise, a test could be too weak to detect the existence of a defect in the generated code. For example, as pointed out by prior work (Liu et al., 2023a), existing benchmarks such as HumanEval (Chen et al., 2021a) and MBPP (Austin et al., 2021) still suffer from “insufficient tests”, allowing incorrect code to pass all the tests without capturing the bugs.

Despite its importance, as shown in Figure 27, among the benchmarks that use test cases as the oracle, only 8.7% considered and reported “test coverage” explicitly in their papers, while 87.8% ignored the test coverage.

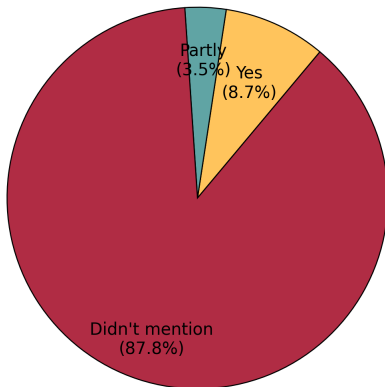


Figure 27: Benchmark Distribution over Test Coverage

Furthermore, we dived into 30 representative benchmarks (Listed in Appendix D) and identified an example (Figure 28) from MBPP (Austin et al., 2021) where the test is incorrect. It alerts us that both the quality of the test and the test adequacy (e.g., code coverage) should be considered.

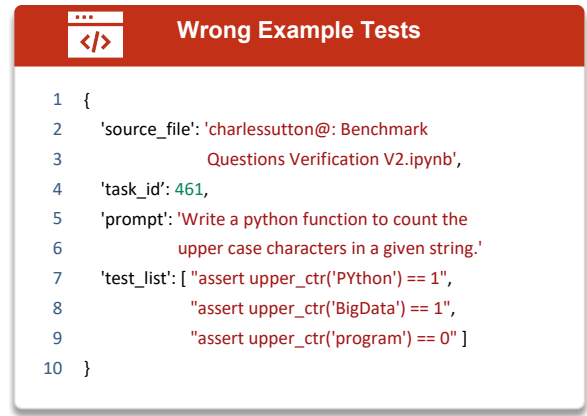


Figure 28: An Example of Incorrect Tests from MBPP (Austin et al., 2021).

B.4 Statistics about Evaluation

Studied LLMs. We summarize the number of LLMs that have been evaluated in each benchmark evaluation. Among the 274 benchmarks, 183 of them are evaluated over LLMs, so we show the statistics over them. As shown in Figure 29, over 34% of the benchmarks were evaluated on fewer than 3 LLMs, with 11.48% benchmarks only evaluated on one LLM. Such evaluation results can hardly be generalized to other LLMs. Furthermore, *more than half of the benchmarks studied fewer than 6 LLMs* ($51\% = (21 + 22 + 20 + 4 + 12 + 15)/183$).

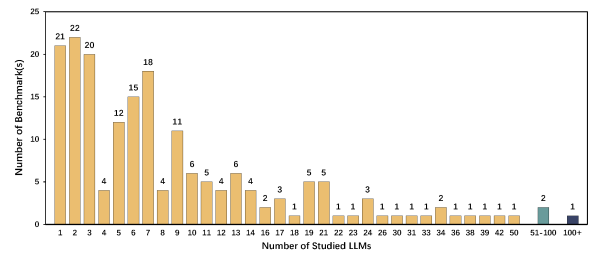


Figure 29: Benchmark Distribution over LLM Experimented

Additionally, we listed the top-10 LLMs by the number of code-related benchmarks they have been evaluated, as shown in Figure 30. GPT series leads significantly with 116 benchmarks, suggesting its widespread adoption and possibly its versatility or superior performance in handling code-related tasks. The rest, including CodeLlama, StarCoder, CodeGen, and others, show varying degrees of involvement, with numbers ranging from 60 down to 24 benchmarks for Claude. This figure may provide a reference for choosing which model to

evaluate. In addition, it is worth mentioning that different LLMs should be considered considering different coding tasks.

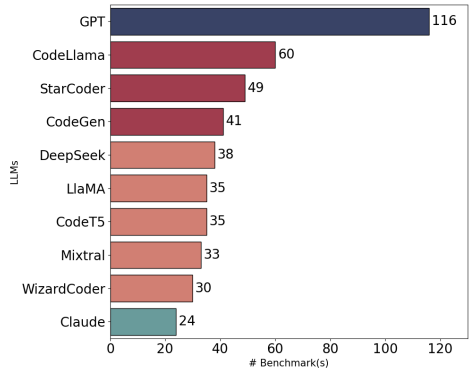


Figure 30: Top-10 Studied LLMs for Code-related Benchmarks

Experiment Environments. The experimental environment (such as the operating system and hardware) is important for the reproduction of the experiment. However, Figure 32 and Figure 31 highlight a significant gap. A mere 27.4% of benchmarks document the devices used in their experiments, leaving a substantial 72.6% that do not. The situation appears even more dire when considering os, with only 3.6% of benchmarks documenting the OS used, while a staggering 96.4% neglect to record this information.

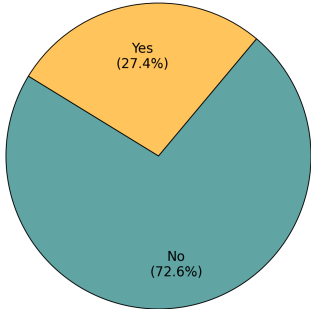


Figure 31: Benchmark Distribution over Recording Experiment Devices

Prompting and Prompting Strategies Prompting has a direct impact on the quality of the LLMs' output results (Wei et al., 2022; He et al., 2024a; Jin et al., 2024). So, we summarized whether different prompting strategies have been evaluated and statistics the distribution. Figure 33 shows the usage of four kinds of prompts: zero-shot, few-shot, chain-of-thought, and retrievals (RAG). From Figure 33, we can see that a vast majority (94.9%) benchmarks were evaluated in a zero-shot context setting, while only 21.2% benchmarks were evalu-

ated in a few-shot manner. Even fewer benchmarks were evaluated under the COT and RAG settings, utilized by only 8.8% and 2.6%.

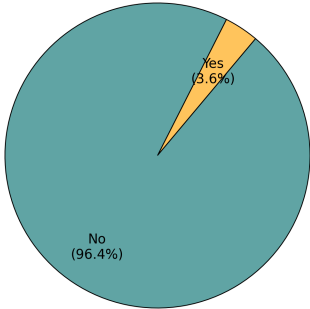


Figure 32: Benchmark Distribution over Recording Experiment OS

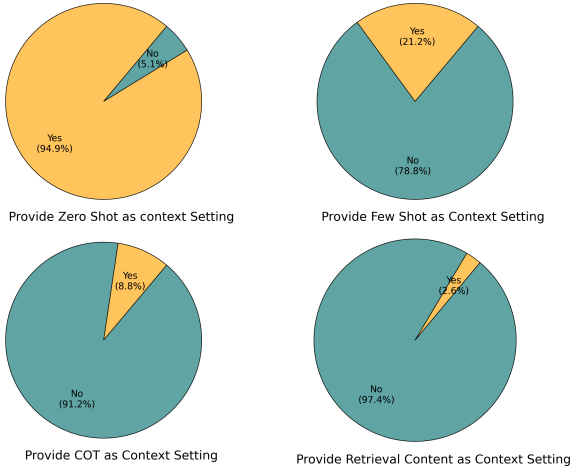


Figure 33: Benchmark Distribution over Context Setting

Prompt Quality The prompt quality also greatly impacts the LLM evaluation (He et al., 2024b). So, carefully designing a prompt needs consideration. However, as shown in Figure 34, 73.3% representative benchmarks (Appendix D) do not validate whether the prompt they used is well-designed.

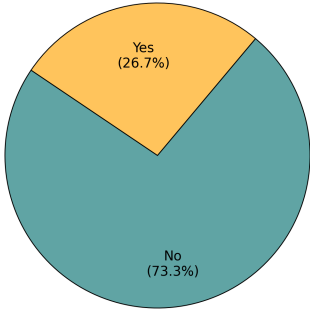


Figure 34: Benchmark Distribution Over Validation of Prompts

Repeated Experiment Given the random na-

ture of LLMs, the experiments are expected to repeat, ensuring the stability and reliability of the results. However, as shown in Figure 35, **only 35.4% benchmarks went through a repeated experiment**, while a majority of 64.6% opted against repeating the experiment. This reflects a lack of awareness regarding the stability and reproducibility of evaluations.

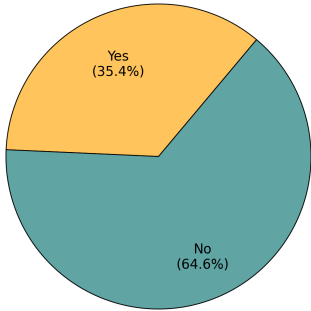


Figure 35: Benchmark Distribution over Repeating the Experiment

B.5 Statistics about Analysis

Experiment Explanation. Explaining experiment results is crucial for other practitioners to understand what the outcomes mean in the context of the research questions. So, we investigate whether the representative benchmarks (Appendix D) have explained the experiment results. As shown in Figure 36, 70% benchmarks have detailed explanations and analyses of their evaluation results, while still 30% have not. Indeed, an explanation contributes to the body of knowledge by making it possible to understand and compare results with previous studies, promoting transparency within the community.

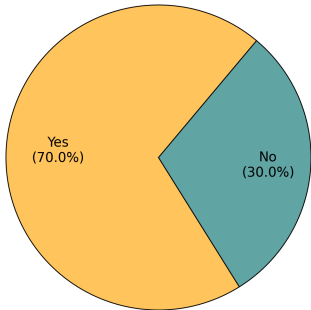


Figure 36: Benchmark Distribution Over Explaining the Experiment

A clear and precise presentation of experimental results is important for enabling robust interpretation and comparison across benchmarks. How-

ever, further examination of the 30 representative benchmarks (listed in Appendix D) revealed notable deficiencies in result visualization. As shown in Figure 37, CruxEval (Gu et al., 2024) exhibits unclear experimental result presentation. Specifically, the scatter plot suffers from ambiguous labeling, poor readability of axis values, and inconsistent marker encoding, making it difficult for researchers to extract meaningful insights. Such presentation shortcomings obscure the performance relationships between methods and compromise the benchmark’s usability for fair evaluation. To address these issues, benchmarks should adopt standardized and well-documented visualization practices, ensuring results are interpretable, accessible, and reproducible.

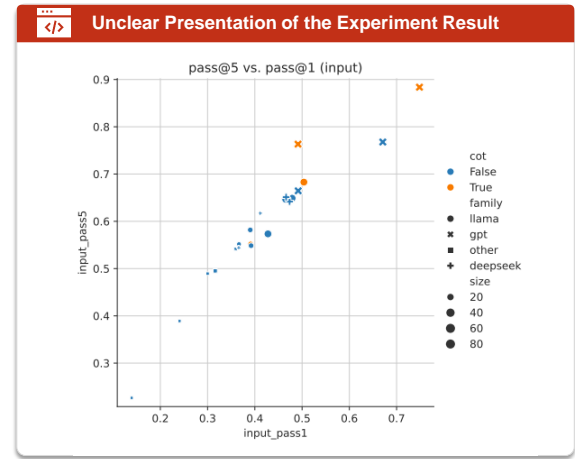


Figure 37: An Example of Unclear Experiment Analysis and Display from CruxEval (Gu et al., 2024)

B.6 Statistics about Release

Data Accessibility. The fundamental requirement for releasing a benchmark is that it must be open-sourced. However, surprisingly, as shown in Figure 38, we observed that 5.1% of the benchmarks are only partially open-sourced (e.g., missing some subjects or tests), and **5.8% are not open-sourced at all** (e.g., links/web pages are no longer active).

Prompt Accessibility. Detailed prompts are essential for ensuring the reproducibility and transparency of code-related benchmarks. Yet, Figure 39 indicates that **52.6% of benchmarks do not provide detailed prompts**, limiting the ability to accurately replicate and evaluate the performance of LLMs. Lacking prompt disclosure highlights a gap in benchmark design practices, as prompts are often indispensable for understanding model perfor-

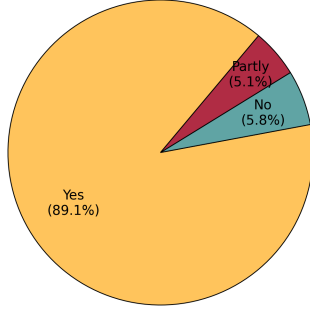


Figure 38: Benchmark Data Availability

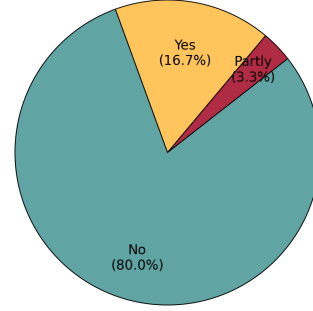


Figure 40: Availability of Logging Information

mance under specific conditions. While 47.4% of benchmarks include such prompts, the absence of comprehensive prompt documentation in over half of the cases raises concerns about the consistency and reproducibility of reported results.

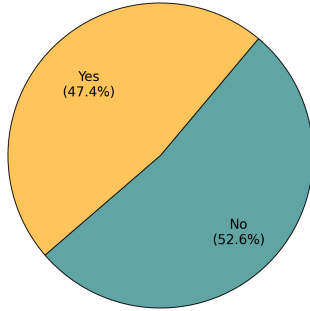


Figure 39: Availability of Prompts

Logging Info Accessibility. Providing detailed logging information, including comprehensive experimental results, is essential for ensuring transparency, verifiability, and reproducibility in benchmarking research. However, as shown in Figure 40, only **16.7% of the benchmarks make their experimental results publicly available**, while 80.0% fail to disclose this critical information. Alarming, an additional 3.3% provide only partial logging details, further complicating result verification. The absence of complete logging information creates significant barriers for researchers attempting to reproduce experiments or validate reported findings, thereby undermining the reliability of benchmarks. To address this, we emphasize the necessity of making detailed logging information, including intermediate results and metrics, publicly accessible to uphold rigorous scientific standards and foster trustworthy comparisons across models.

User Manual Accessibility. A high-quality user manual, such as a well-documented README file, is crucial for enhancing benchmark usability, enabling users to understand the dataset, ex-

ecute provided scripts, and reproduce results efficiently. However, our analysis revealed that a significant number of benchmarks lack comprehensive user manuals, hindering accessibility and adoption. As depicted in Figure 41, poorly structured or incomplete manuals often omit essential components such as benchmark introductions, usage instructions, and evaluation scripts. This creates unnecessary barriers for researchers who rely on these manuals for setup and experimentation. To address this, we advocate for benchmarks to include clear, standardized user manuals that provide an overview of the benchmark, step-by-step execution guides, and troubleshooting instructions, ensuring a seamless and reproducible user experience.

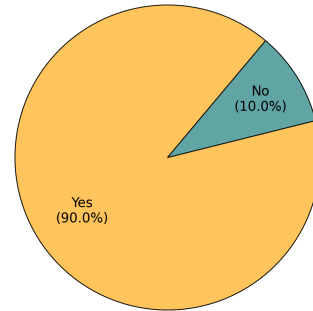


Figure 41: Availability of User Manual

Convenient Evaluation Interface Availability. Providing convenient evaluation interfaces is essential for enhancing the usability and accessibility of benchmarks, enabling researchers to easily reproduce results and compare models. As shown in Figure 42, **16.7% of benchmarks fail to offer any evaluation interfaces**, imposing significant barriers to usability. While a majority of benchmarks (83.3%) provide such interfaces, including command-line tools, Docker images, or scripts, the absence of standardized and user-friendly evaluation tools in a notable minority of cases highlights an area for improvement. Benchmarks without

convenient evaluation interfaces require users to spend additional effort in setup and result verification, which can discourage adoption and hinder reproducibility. To address this, we emphasize the importance of releasing benchmarks with well-documented, ready-to-use evaluation pipelines to promote efficient, reliable, and fair benchmarking practices.

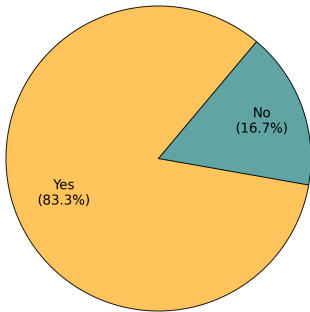


Figure 42: Availability of Convenient Evaluation Interfaces

Temperature Records. One critical parameter for benchmarking is the temperature setting, which influences stochasticity in LLMs. As shown in Figure 43, we observed that **57.3% of benchmarks fail to record the temperature setting**, hindering reproducibility and fair evaluation. While 42.7% of benchmarks do document this parameter, the majority omission highlights an overlooked yet essential aspect of benchmark transparency.

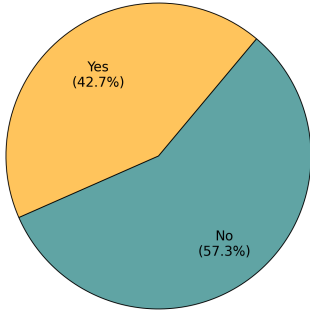


Figure 43: Benchmark Distribution over Recording Temperature

License Provision. Releasing benchmarks under a clear and accessible license is fundamental for legal compliance and ensuring open collaboration. Figure 44 reveals that **19.3% of benchmarks do not provide a license**, limiting their usability and distribution. Encouragingly, 80.7% of benchmarks do include a license, but the lack of licensing in nearly one-fifth of the benchmarks raises concerns about widespread adoption and usage. These find-

ings emphasize the need for standardized practices in benchmark releases to promote legal clarity and accessibility.

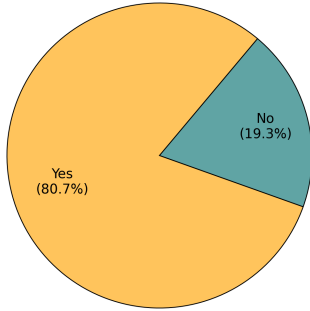


Figure 44: Provision of License

Data Security. Ensuring data security is a critical yet often overlooked aspect of benchmark development. Sensitive information, such as API keys, credentials, or private tokens, should never be included in benchmark releases. However, further investigation into 30 representative benchmarks (listed in Appendix D) revealed instances of sensitive data leakage. As shown in Figure 45, XSemPLR (Zhang et al., 2023b) inadvertently included an **API key** in its release, a critical oversight that can expose resources to external exploitation. Similarly, Figure 46 highlights an example from

```
API Key Leakage
1  #!/bin/bash
2
3  # conda activate skg
4  #export WANDB_API_KEY=*****
5  export WANDB_PROJECT=mt5-large_mgequery_few-shot
6  export CUDA_LAUNCH_BLOCKING=1
```

Figure 45: An Example of API Key Leakage in Benchmark Release from XSemPLR (Zhang et al., 2023b).

CrossVul (Nikitopoulos et al., 2021), where **personal names and email addresses** were unintentionally disclosed. Such leakage poses risks of unauthorized access and resource misuse, potentially compromising systems and research integrity.

Usability. Clear and comprehensive documentation is crucial for ensuring the usability of benchmarks, as poorly written instructions can significantly hinder adoption and reproducibility. We dived into the 30 representative benchmarks (listed in Appendix D) and identified an example where

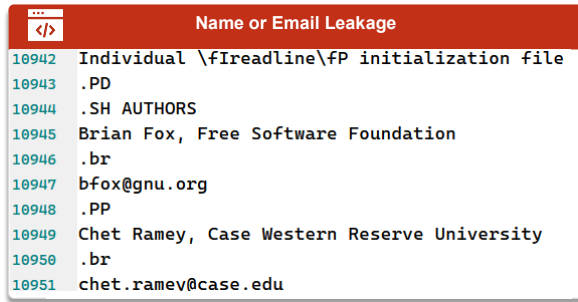


Figure 46: An Example of Name & Email Leakage in Benchmark Release from CrossVul (Nikitopoulos et al., 2021).

the README file provided insufficient and unclear information. As shown in Figure 47, VulDeePecker (Li et al., 2018b) includes a less-than-ideal ReadMe file, which lacks essential usage instructions and evaluation guidelines, making the benchmark difficult to understand and deploy.

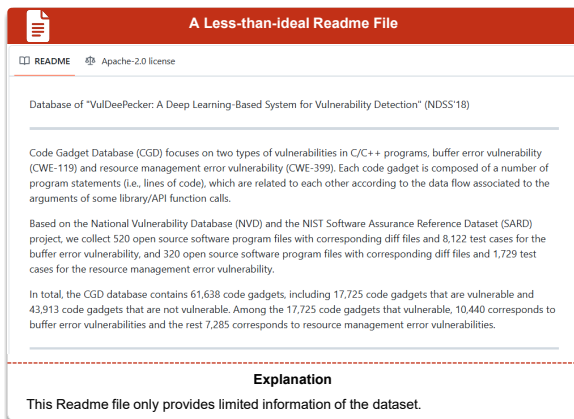


Figure 47: An Example of Unreadable and Hard-to-Use README in Benchmark Release from VulDeePecker (Li et al., 2018b).

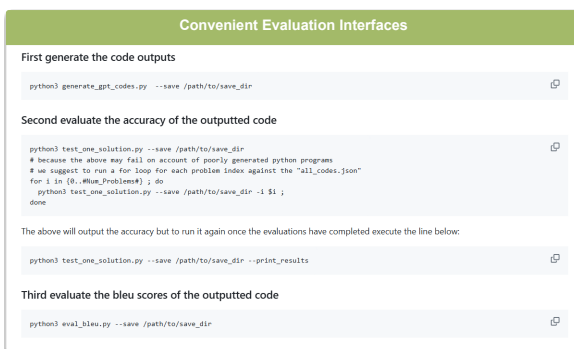


Figure 48: A Good Example of Easy-to-Read README in Benchmark Release from APPS (Hendrycks et al., 2021).

In contrast, Figure 48 highlights APPS (Hendrycks et al., 2021), which provides well-structured and easy-to-follow documentation. The APPS benchmark includes step-by-step instructions for generating, evaluating, and analyzing results, enabling users to efficiently reproduce experiments.

These observations emphasize the importance of high-quality documentation for benchmarks to enhance accessibility, reduce friction in usage, and foster reproducible research.

C Details of Human Study

C.1 Interviewee Selection

The selection of interviewees is pivotal to ensuring the representativeness and relevance of the data collected. This involves identifying individuals with the knowledge or experience pertinent to the research theme.

To this end, we chose graduate students from SE or AI fields who have published at least one paper. This criterion ensures that participants have research experience and judgment capabilities. The focus on SE and AI fields is due to their likely interest in code benchmarks. Particularly, we aimed to recruit individuals who have published papers on code benchmarks to obtain firsthand feedback from experienced benchmark developers.

C.2 Survey Question Design

Questions. The body of the survey was divided into five stages of benchmark development (following Figure 1), with necessary background information provided for each stage. Each criterion in HOW2BENCH was slightly modified to be in the first-person perspective, making it easier for interviewees to empathize and answer the questions from their own viewpoint. Finally, to facilitate comprehension, questions and instructions were translated into both English and Chinese.

Question Setting. To minimize the effort required from respondents, we designed *single-choice questions* with four options:

☐ I found it **important**, and I **have done** it.

☐ I found it **important**, although I **haven't done** it.

☐ I found it **not important**, but I **have done** it.

☐ I found it **not important**, and I **wouldn't** do it.

This format is intended to orthogonally explore the correlation between *awareness* and *behavior*.

C.3 Interview Process

Questionnaire Distribution The questionnaire was distributed via online platforms, targeting academic and professional networks related to SE and AI. The distribution started on October 27, 2024, and ended on November 27th, 2024, lasting one month.

Results Collection The responses were automatically collected through the online platform used for distribution.

Survey Screening Since the requirement was for participants who have published papers, responses

Geographical distribution of Interviewees

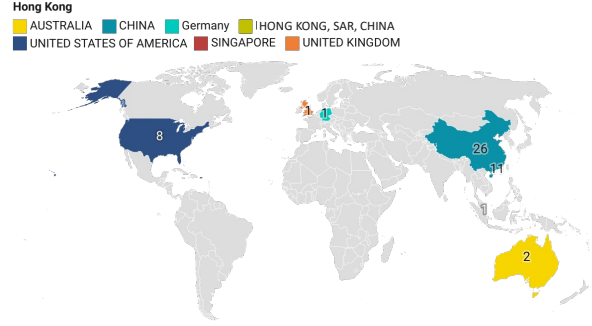


Figure 49: Geographical Distribution of Interviewees

from those selecting “No” to having published a paper were excluded. Also, incomplete surveys where not all questions were answered were also considered invalid and excluded from the analysis.

C.4 Interview Result Analysis

In total, we collected 50 responses. The respondents were from seven regions, including the United States, the United Kingdom, Germany, Australia, China, and others, as shown in Figure 49. Only one survey was invalid due to the respondent selecting “have not published a paper”, leaving **49 valid surveys** for analysis. A breakdown of the respondents’ demographics is shown in Figure 50. The detailed responses for all 55 criteria in HOW2BENCH are shown in Figure 51 and Figure 52.

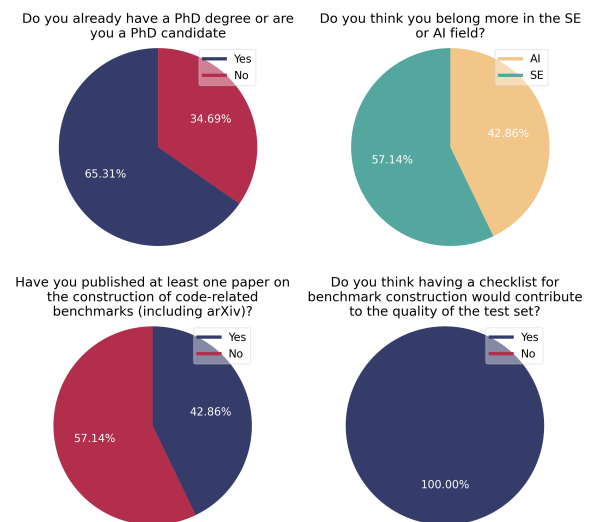


Figure 50: Demography of Interviewees

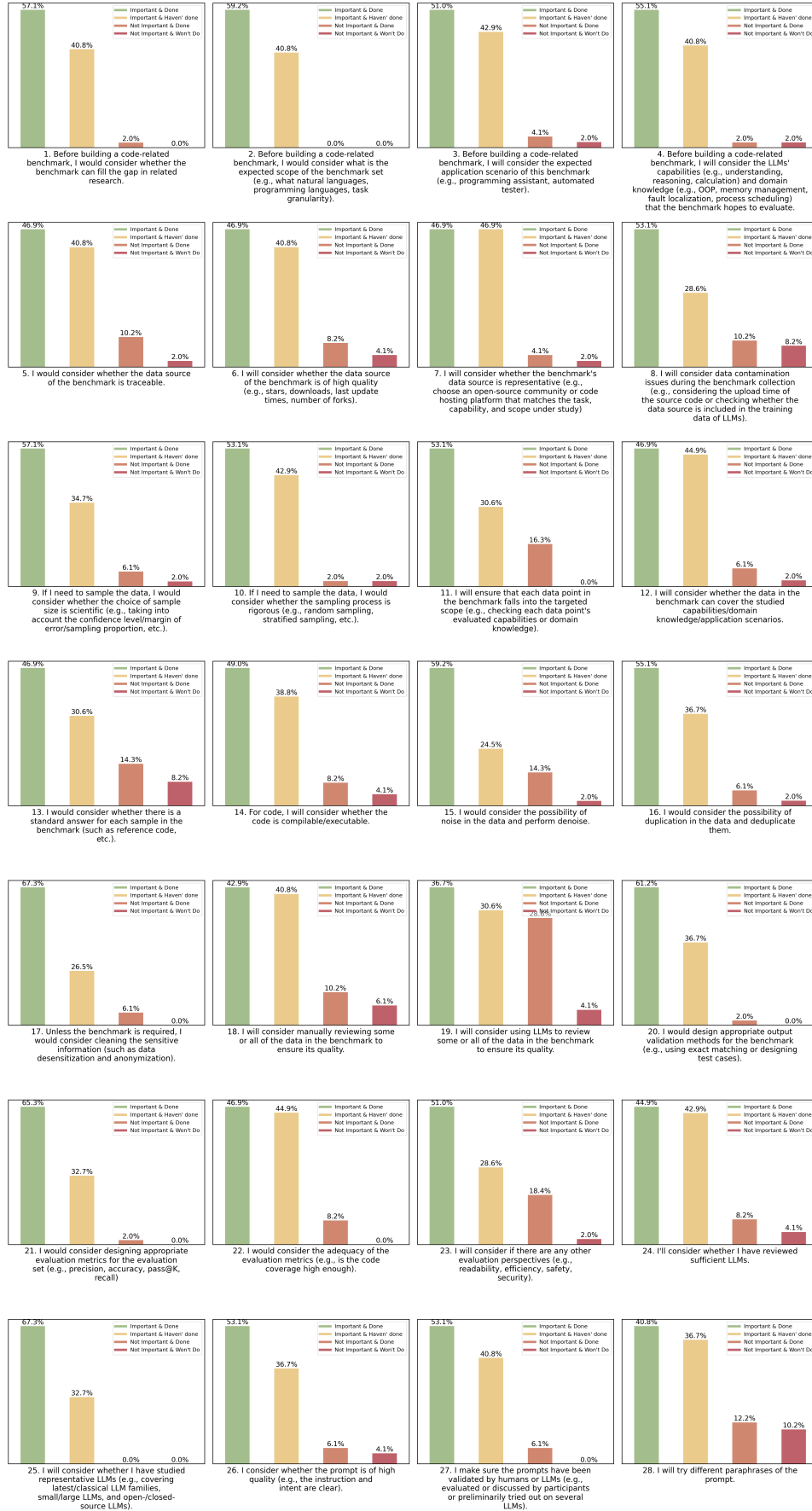


Figure 51: Results of Human Study (Questions 1 - 28)

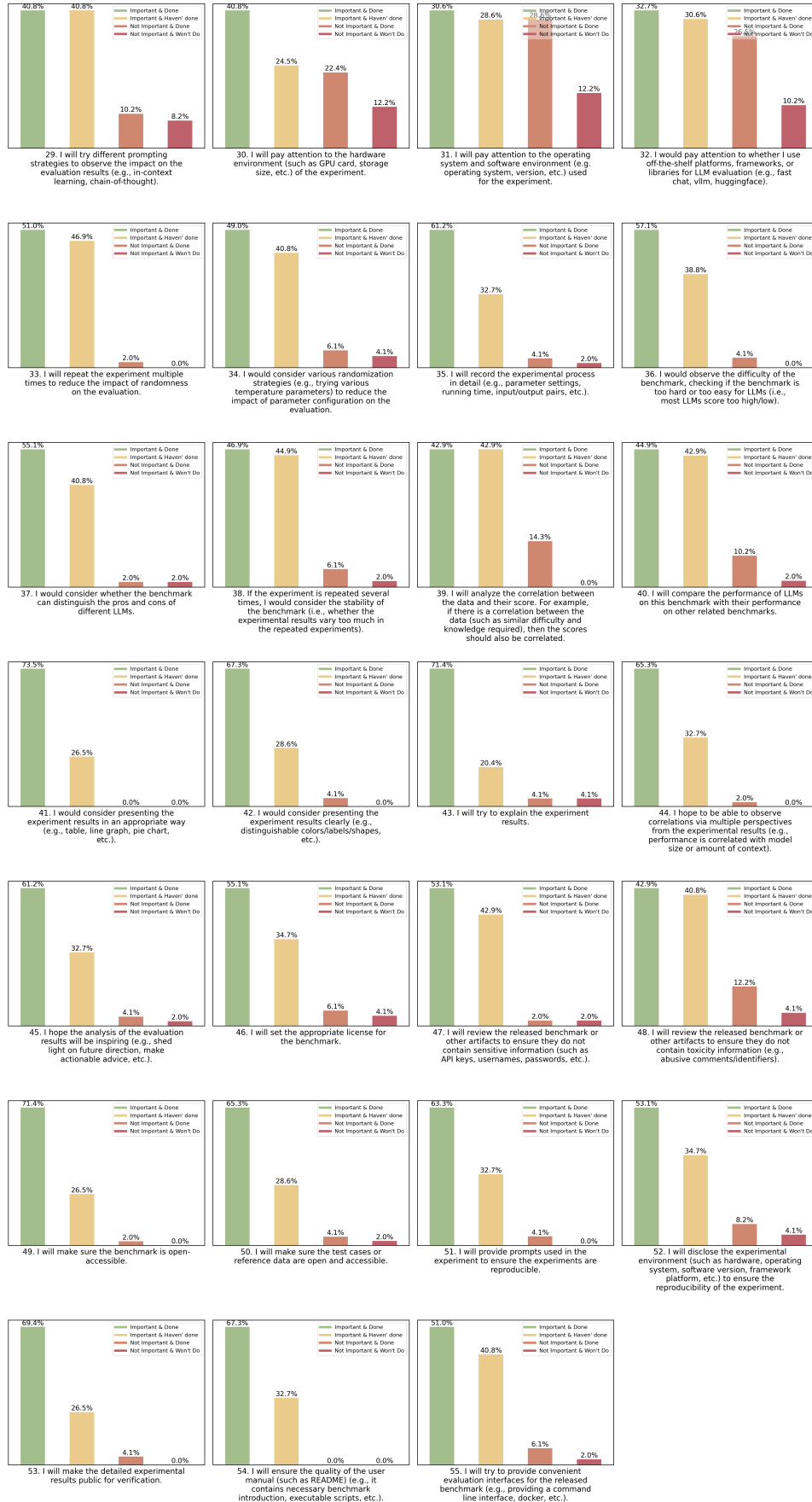


Figure 52: Results of Human Study (Questions 29 - 55)

2792	D List of Studied Benchmarks (Focused	The latest one as of 31/8/2024:	2825
2793	Ones)		
2794	Code Generation: Five with top citations:	• Long Code Arena (Bogomolov et al., 2024)	2826
2795	• HumanEval (Chen et al., 2021a)	Text To SQL: Five with top citations:	2827
2796	• MBPP (Austin et al., 2021)	• WikiSQL (Zhong et al., 2017)	2828
2797	• CodeContest (Li et al., 2022)	• Spider (Yu et al., 2018)	2829
2798	• leetcodehardgym (Shinn et al., 2023)	• Advising (Finegan-Dollak et al., 2018)	2830
2799	• APPS (Hendrycks et al., 2021)	• BIRD (Li et al., 2023a)	2831
2800	The latest one as of 31/8/2024:	• Spider-Realistic (Deng et al., 2021)	2832
2801	• VerilogEval (Pinckney et al., 2024)	The latest one as of 31/8/2024:	2833
2802	Defect Detection: Five with top citations:	• AMBROSIA (Saparina and Lapata, 2024)	2834
2803	• VulDeePecker (Li et al., 2018b)		
2804	• Devign (Zhou et al., 2019)		
2805	• Chromium and Debian (Chakraborty et al., 2022)		
2806	• μ VulDeePecker (Zou et al., 2020)		
2807	• Synthetic Dataset (Hellendoorn et al., 2020)		
2808	The latest one as of 31/8/2024:		
2809	• VulDetectBench (Liu et al., 2024c)		
2810	Program Repair: Five with top citations:		
2811	• Defects4J (Just et al., 2014)		
2812	• BFP (Tufano et al., 2019)		
2813	• MANYBUGS, INTROCLASS (Le Goues et al.,		
2814	2015)		
2815	• HumanEval-Java (Jiang et al., 2023)		
2816	• QuixBugs (Prenner et al., 2022)		
2817	The latest one as of 31/8/2024:		
2818	• DebugBench (Tian et al., 2024)		
2819	Code Summarization: Five with top citations:		
2820	• CODE-NN (Iyer et al., 2016)		
2821	• Java-small/med/large (Alon et al., 2019)		
2822	• code-summarization-public (Wan et al., 2018)		
2823	• HumanEvalPack (Muennighoff et al., 2024)		
2824	• Shrivastava et al. (Shrivastava et al., 2023b)		

E List of Studied Benchmarks (Full)

We collected and studied 274 code-related benchmarks. We then listed and grouped them by year. **2024:**

- CodeEditorBench (Guo et al., 2024)
- MHPP (Dai et al., 2024)
- LiveCodeBench (Jain et al., 2024)
- CodeAgentBench (Zhang et al., 2024a)
- CruxEval (Gu et al., 2024)
- BigCodeBench (Zhuo et al., 2024)
- OOPEval (Wang et al., 2024b)
- DevEval (Li et al., 2024a)
- Long Code Arena (Bogomolov et al., 2024)
- CodeRAGBench (Wang et al., 2024d)
- ScenEval (Paul et al., 2024)
- AICoderEval (Xia et al., 2024b)
- VersiCode (Wu et al., 2024b)
- VHDL-Eval (Vijayaraghavan et al., 2024)
- NaturalCodeBench (Zhang et al., 2024b)
- CodeGuard+ (Fu et al., 2024)
- PECC (Haller et al., 2024)
- USACO (Shi et al., 2024b)
- ParEval (Nichols et al., 2024)
- MxEval (Athiwaratkun et al., 2022)
- MMCode (Li et al., 2024c)
- Plot2Code (Wu et al., 2024a)
- ChartMimic (Shi et al., 2024a)
- DebugBench (Tian et al., 2024)
- PythonIO (Zhang et al., 2024c)
- StaCCQA (Yang et al., 2024a)
- RepoQA (Liu et al., 2024a)
- PRIMEVUL (Ding et al., 2024a)
- VulDetectBench (Liu et al., 2024c)

- ProCQA (Li et al., 2024e)
- CoSQA+ (Gong et al., 2024) (Huang et al., 2021)
- JavaBench (Cao et al., 2024a)
- HumanEvo (Zheng et al., 2024)
- REPOEXEC (Hai et al., 2024)
- EHR-SeqSQL (Ryu et al., 2024)
- BookSQL (Kumar et al., 2024)
- AMBROSIA (Saparina and Lapata, 2024)
- WUB, WCGB (Yun et al., 2024)
- RES-Q (LaBash et al., 2024)
- PythonSaga (Yadav et al., 2024b)
- Mercury (Du et al., 2024)
- ENAMEL (Qiu et al., 2024b)
- RealHumanEval (Mozannar et al., 2024)
- CoderUJB (Zeng et al., 2024)
- EvoEval (Xia et al., 2024a)
- ML-Bench (Liu et al., 2023c)
- VerilogEval (Pinckney et al., 2024)
- CodeApex (Fu et al., 2023)
- HumanEvalPack (Muennighoff et al., 2024)
- HumanEval+ (Liu et al., 2023b)
- HumanEval-X (Zheng et al., 2023a)
- XCodeEval (Khan et al., 2024)
- CoderEval (Yu et al., 2023)
- CodeXGLUE (Lu et al., 2021)
- VulnPatchPairs (Risse and Böhme, 2024)
- WikiSQL (Zhong et al., 2017)
- CrossCodeEval (Ding et al., 2023)
- SWE-bench (Jimenez et al., 2024)
- BAIRI et al. (Bairi et al., 2024)
- BioCoder (Tang et al., 2024)
- RepoBench (Liu et al., 2024b)

2900	• NoFunEval (Singhal et al., 2024)	• CodeTransOcean (Yan et al., 2023)	2932
2901	• CoCoMIC (Ding et al., 2024b)	• G-TransEval (Jiao et al., 2023)	2933
2902	• Java-small/med/large (Alon et al., 2019)	• AVATAR (Ahmad et al., 2023)	2934
2903	• FixEval (Haque et al., 2023)	• RunBugRun (Prenner and Robbes, 2023)	2935
2904	• CommitBench (Schall et al., 2024)	• VulBench (Gao et al., 2023b)	2936
2905	• InfiAgent-DABench (Hu et al., 2024)	• DiverseVul (Chen et al., 2023)	2937
2906	• InfiBench (Li et al., 2024d)	• Hellendoorn et al. (Hellendoorn et al., 2020)	2938
2907	• Design2Code (Si et al., 2024)	• XSemPLR (Zhang et al., 2023b)	2939
2908	• MatPlotBench (Yang et al., 2024b)	• BIRD (Li et al., 2023a)	2940
2909	• EditEval (Li et al., 2024b)	• Stack-Repo (Shrivastava et al., 2023a)	2941
2910	• D1, D2, D3 (Huang et al., 2024)	• RepoEval (Liao et al., 2024)	2942
2911	• RepoEval (Liao et al., 2024)	• MTPB (Nijkamp et al., 2022)	2943
2912	• BetterTypes4Py, InferTypes4Py (Wei et al., 2023)	• ARCADE (Yin et al., 2023)	2944
2913	• HumanEval-Java (Jiang et al., 2023)	• Shrivastava et al. (Shrivastava et al., 2023b)	2945
2914	• PIE (Shypula et al., 2024)	• Grag et al. (Garg et al., 2022)	2946
2915	• EvalGPTFix (Zhang et al., 2023a)	• GSM-HARD (Gao et al., 2023a)	2947
2916	• EHRSQL (Lee et al., 2023)	• InferredBugs (Jin et al., 2023)	2948
2917	• Spider2-V (Cao et al., 2024c)	• LeetcodeHardGym (Shinn et al., 2023)	2949
2918	• TESTEVAL (Wang et al., 2024c)	• APIBench (Patil et al., 2023)	2950
2919	• ChatTester (Yuan et al., 2023b)	• ClassEval (Du et al., 2023b)	2951
2920	• Code Lingua (Pan et al., 2024)	• CommitChronicle (Eliseeva et al., 2023)	2952
2921	• EffiBench (HUANG et al., 2024)	• TeCo (Nie et al., 2023)	2953
2922	• CRUXEval-X (Xu et al., 2024)	• TESTPILOT (Schäfer et al., 2024)	2954
2923	• DomainEval (Zhu et al., 2024)	2022:	2955
2924	2023:	• AixBench (Hao et al., 2022)	2956
2925	• MCoNaLa (Wang et al., 2023b)	• TypeBugs (Oh and Oh, 2022)	2957
2926	• MultiPL-E (Cassano et al., 2022)	• XLCoST (Zhu et al., 2022)	2958
2927	• ODEX (Wang et al., 2022)	• CS1QA (Lee et al., 2022)	2959
2928	• TACO (Li et al., 2023b)	• Chromium and Debian (Chakraborty et al., 2022)	2960
2929	• DOTPROMPTS,	• Spider-Realistic (Deng et al., 2021)	2961
2930	CROBENCH (Agrawal et al., 2023)	• Spider-SS (Gan et al., 2022)	2962
2931	• StudentEval (Babe et al., 2024)	• DSP (Chandel et al., 2022)	2963

2964	• CodeContest (Li et al., 2022)	• PACS (Heyman and Cutsem, 2020)	2997
2965	• PandasEval, NumpyEval (Zan et al., 2022b)	• Criteria2SQL (Yu et al., 2020)	2998
2966	• TorchDataEval, MonkeyEval, BeatNu-	• SQUALL (Shi et al., 2020)	2999
2967	mEval (Zan et al., 2022a)	• Hu et al. (Hu et al., 2019)	3000
2968	• DS-1000 (Lai et al., 2023)	• CodeSearchNet Challenge (Husain et al., 2019)	3001
2969	• MCMD (Tao et al., 2022)	• MIMICSQL (Wang et al., 2020)	3002
2970	• ExeDS (Huang et al., 2022)	• Atlas (Watson et al., 2020)	3003
2971	• QuixBugs (Prenner et al., 2022)	• Liu et al. (Liu et al., 2022)	3004
2972	• ManyTypes4Py v0.7 (Mir et al., 2022)	• Android (Agarwal et al., 2020)	3005
2973	2021:	• CCSD (Liu et al., 2021)	3006
2974	• SySeVR (Li et al., 2018a)	2019:	3007
2975	• Ling&Wu et al. (Ling et al., 2021)	• BFP (Tufano et al., 2019)	3008
2976	• Chen et al. (Chen et al., 2021b)	• SARD (Lin et al., 2019)	3009
2977	• MBPP, MathQA-Python (Austin et al., 2021)	• Spider (Yu et al., 2018)	3010
2978	• HumanEval (Chen et al., 2021a)	• JuICe (Agashe et al., 2019)	3011
2979	• APPS (Hendrycks et al., 2021)	• Nguyen et al. (Nguyen et al., 2019)	3012
2980	• Berabi et al. (Berabi et al., 2021)	• Lin et al. (Lin et al., 2021)	3013
2981	• CrossVul (Nikitopoulos et al., 2021)	• Zhou et al. (Zhou et al., 2019)	3014
2982	• PYPIBUGS, RANDOMBUGS (Allamanis et al.,	• CoSQL (Yu et al., 2019a)	3015
2983	2021)	• SParC (Yu et al., 2019b)	3016
2984	• D2A (Zheng et al., 2021)	• Malik (Malik et al., 2019)	3017
2985	• CodeQA (Liu and Wan, 2021)	• LeClair (LeClair et al., 2019)	3018
2986	• Spider-DK (Gan et al., 2021b)	2018:	3019
2987	• KaggleDBQA (Lee et al., 2021)	• CoNaLa (Yin et al., 2018)	3020
2988	• SEDE (Hazoom et al., 2021)	• DeepCom (Hu et al., 2018a)	3021
2989	• Spider-Syn (Gan et al., 2021a)	• TL-CodeSum (Hu et al., 2018b)	3022
2990	• CoDesc (Hasan et al., 2021)	• code-summarization-public (Wan et al., 2018)	3023
2991	• Methods2Test (Tufano et al., 2022)	• Russell et al. (Russell et al., 2018)	3024
2992	• Rozière et al. (Rozière et al., 2022)	• VulDeePecker (Li et al., 2018b)	3025
2993	2020:	• Lin et al. (Lin et al., 2018)	3026
2994	• Lachaux&Roziere et al. (Rozière et al., 2020)	• StaQC (Yao et al., 2018)	3027
2995	• μ VulDeePecker (Zou et al., 2020)	• Advising (Finegan-Dollak et al., 2018)	3028
2996	• CosBench (Yan et al., 2020)		

• ConCode (Iyer et al., 2018)

• NNGen (Liu et al., 2018)

• Gu et al. (Gu et al., 2018)

2017:

• QuixBugs (Lin et al., 2017)

• the DeepFix dataset (Gupta et al., 2017)

• Barone et al. (Barone and Sennrich, 2017)

2016:

• CODE-NN (Iyer et al., 2016)

• Mou et al. (Mou et al., 2016)

2015:

• MANYBUGS, INTROCLASS (Le Goues et al., 2015)

2014:

• Defects4j (Just et al., 2014)

• BigCloneBench (Svajlenko et al., 2014)

F Guideline

Finally, for ease of printing and use, we organized the guideline HOW2BENCH into a clear, color-coded checklist (4 pages in total) that is easy to print, attached at the end of the paper.

HOW-TO-BENCH (1/4)			
	Phase 1. Benchmark Design	Priority	<input checked="" type="checkbox"/>
1	Consider whether the benchmark can <u>fill the gap in related research</u> .	★★★	<input type="checkbox"/>
2	Consider what is the <u>expected scope</u> of the benchmark set (e.g., what natural languages, programming languages, task granularity).	★★★	<input type="checkbox"/>
3	Consider the <u>expected application scenario</u> of this benchmark (e.g., programming assistant, automated tester).	★★★	<input type="checkbox"/>
4	Consider the <u>LLMs' capabilities</u> (e.g., understanding, reasoning, calculation) and <u>domain knowledge</u> (e.g., OOP, memory management, fault localization, process scheduling) that the benchmark hopes to evaluate.	★★★	<input type="checkbox"/>
	Phase 2. Benchmark Construction (1/2)	Priority	<input checked="" type="checkbox"/>
5	Consider whether the <u>data source</u> of the benchmark is <u>traceable</u> .	★	<input type="checkbox"/>
6	Consider whether the data source of the benchmark is of <u>high quality</u> (e.g., stars, downloads, last update times, number of forks).	★★	<input type="checkbox"/>
7	Consider whether the benchmark's data source is <u>representative</u> (e.g., choose an open-source community or code hosting platform that matches the task, capability, and scope under study)	★	<input type="checkbox"/>
8	Consider <u>data contamination issues</u> during the benchmark collection (e.g., considering the upload time of the source code or checking whether the data source is included in the training data of LLMs).	★	<input type="checkbox"/>
9	If data <u>sampling</u> is needed, consider whether the <u>choice of sample size is scientific</u> (e.g., considering the confidence level/margin of error/sampling proportion, etc.).	★	<input type="checkbox"/>
10	If data <u>sampling</u> is needed, consider whether the <u>sampling process is rigorous</u> (e.g., random sampling, stratified sampling, etc.).	★	<input type="checkbox"/>
11	Ensure each data point in the benchmark <u>falls into the targeted scope</u> (e.g., checking each data point's evaluated capabilities or domain knowledge).	★	<input type="checkbox"/>
12	Consider whether the data in the benchmark can <u>cover</u> the studied capabilities/domain knowledge/application scenarios.	★★★	<input type="checkbox"/>
13	Consider whether there is a <u>standard answer</u> for each sample in the benchmark (such as reference code, etc.).	★★	<input type="checkbox"/>
14	For <u>code</u> , consider whether the code is <u>compilable/executable</u> .	★	<input type="checkbox"/>
15	Consider the possibility of noise in the data and perform <u>denoise</u> .	★★	<input type="checkbox"/>
16	Consider the possibility of duplication in the data and <u>deduplicate</u> them.	★★	<input type="checkbox"/>
17	<u>Clean the sensitive information</u> (such as data desensitization and anonymization) unless the benchmark is deliberately designed so.	★★	<input type="checkbox"/>

HOW-TO-BENCH (2/4)			
	Phase 2. Benchmark Construction (2 / 2)	Priority	<input checked="" type="checkbox"/>
18	<u>Manually review</u> some or all of the data in the benchmark to ensure its quality.	★★	<input type="checkbox"/>
19	<u>Use LLMs to review</u> some or all of the data in the benchmark to ensure its quality.	★	<input type="checkbox"/>
20	Design appropriate <u>output validation methods</u> for the benchmark (e.g., using exact matching or designing test cases).	★★★	<input type="checkbox"/>
21	Design appropriate <u>evaluation metrics</u> for the evaluation set (e.g., precision, accuracy, pass@K, recall).	★★★	<input type="checkbox"/>
22	Consider <u>the adequacy of the evaluation metrics</u> (e.g., is the code coverage high enough).	★★★	<input type="checkbox"/>
23	Consider if there are any <u>other evaluation perspectives</u> (e.g., readability, efficiency, safety, security).	★	<input type="checkbox"/>
	Phase 3. Benchmark Evaluation	Priority	<input checked="" type="checkbox"/>
24	Consider whether <u>sufficient</u> LLMs are evaluated.	★	<input type="checkbox"/>
25	Consider whether <u>representative</u> LLMs (e.g., covering latest/classical LLM families, small/large LLMs, and open-/closed-source LLMs) are evaluated.	★	<input type="checkbox"/>
26	Consider whether the prompt is of <u>high quality</u> (e.g., the instruction and intent are clear).	★★★	<input type="checkbox"/>
27	The prompts have been <u>validated by humans or LLMs</u> (e.g., evaluated or discussed by participants or preliminarily tried out on several LLMs).	★	<input type="checkbox"/>
28	Try <u>different paraphrases</u> of the prompt.	★	<input type="checkbox"/>
29	Try <u>different prompting strategies</u> to observe the impact on the evaluation results (e.g., in-context learning, chain-of-thought).	★★	<input type="checkbox"/>
30	Pay attention to the <u>hardware environment</u> (such as GPU card, storage size, etc.) of the experiment.	★★★	<input type="checkbox"/>
31	Pay attention to the <u>operating system and software environment</u> (e.g. operating system, version, etc.) used for the experiment.	★★★	<input type="checkbox"/>
32	Pay attention to the off-the-shelf <u>platforms, frameworks, or libraries</u> for LLM evaluation (e.g., fast chat, vllm, huggingface) that are used.	★★	<input type="checkbox"/>
33	<u>Repeat</u> the experiment multiple times to reduce the impact of <u>randomness</u> on the evaluation.	★	<input type="checkbox"/>
34	Consider various <u>randomization strategies</u> (e.g., trying various temperature parameters) to reduce the impact of parameter configuration on the evaluation.	★★	<input type="checkbox"/>
35	<u>Record</u> the experimental process in detail (e.g., parameter settings, running time, input/output pairs, etc.).	★★★	<input type="checkbox"/>

HOW-TO-BENCH (3/4)			
	Phase 4. Benchmark Analysis	Priority	<input checked="" type="checkbox"/>
36	Observe the <u>difficulty</u> of the benchmark, checking if the benchmark is too hard or too easy for LLMs (i.e., most LLMs score too high/low).	★★	<input type="checkbox"/>
37	Consider whether the benchmark can <u>distinguish</u> the pros and cons of different LLMs.	★	<input type="checkbox"/>
38	If the experiment is <u>repeated several times</u> , consider the <u>stability of the benchmark</u> (i.e., whether the experimental results vary too much in the repeated experiments).	★	<input type="checkbox"/>
39	Analyze <u>the correlation between the data and their score</u> . For example, if there is a correlation between the data (such as similar difficulty and knowledge required), then the scores should also be correlated.	★★	<input type="checkbox"/>
40	Compare the performance of LLMs on this benchmark with their performance on other related benchmarks .	★	<input type="checkbox"/>
41	Consider <u>presenting the experiment results in an appropriate way</u> (e.g., table, line graph, pie chart, etc.).	★★★	<input type="checkbox"/>
42	Consider <u>presenting the experiment results clearly</u> (e.g., distinguishable colors/labels/shapes, etc.).	★★★	<input type="checkbox"/>
43	<u>Explain</u> the experiment results.	★★★	<input type="checkbox"/>
44	Observe <u>correlations via multiple perspectives</u> from the experimental results (e.g., performance is correlated with model size or amount of context).	★	<input type="checkbox"/>

HOW-TO-BENCH (4/4)			
	Phase 5. Benchmark Release	Priority	<input checked="" type="checkbox"/>
45	The analysis of the evaluation results will be <u>inspiring</u> (e.g., shed light on future direction, make actionable advice, etc.).	★	<input type="checkbox"/>
46	Set the appropriate <u>license</u> for the benchmark.	★★★	<input type="checkbox"/>
47	Review the released benchmark or other artifacts to ensure they <u>do NOT contain sensitive information</u> (e.g., API keys, usernames, passwords, etc.).	★★★	<input type="checkbox"/>
48	review the released benchmark or other artifacts to ensure they <u>do NOT contain toxicity information</u> (e.g., abusive comments/identifiers).	★★★	<input type="checkbox"/>
49	Make sure the benchmark is <u>open-accessible</u> .	★★★	<input type="checkbox"/>
50	Make sure the <u>test cases</u> or <u>reference data</u> are open and accessible.	★★★	<input type="checkbox"/>
51	<u>Provide prompts</u> used in the experiment to ensure the experiments are reproducible.	★★★	<input type="checkbox"/>
52	<u>Disclose the experimental environment</u> (e.g., hardware, operating system, software version, framework platform) to ensure the reproducibility of the experiment.	★★★	<input type="checkbox"/>
53	Make the <u>detailed</u> experimental results <u>public</u> for verification.	★★★	<input type="checkbox"/>
54	Ensure the <u>quality</u> of the user manual such as README (e.g., it contains necessary benchmark introduction, executable scripts, etc.).	★★	<input type="checkbox"/>
55	<u>Provide convenient evaluation interfaces</u> for the released benchmark (e.g., providing a command line interface, docker, etc.).	★★	<input type="checkbox"/>