# DevBench: A Comprehensive Benchmark for Software Development

Bowen Li [* 1]   Wenhan Wu[† * 2]   Ziwei Tang[† * 3]   Lin Shi[† * 4]   John Yang[5]   Jinyang Li[6]   Shunyu Yao[5]   Chen Qian[7]
Binyuan Hui   Qicheng Zhang   Zhiyin Yu   He Du   Ping Yang[8]   Dahua Lin[1]   Chao Peng[8]   Kai Chen[1]

## Abstract

Recent advancements in large language models (LLMs) have significantly enhanced their coding capabilities. However, existing benchmarks predominantly focused on simplified or isolated aspects of programming, such as single-file code generation or repository issue debugging, falling short of measuring the full spectrum of challenges raised by real-world programming activities. To this end, we propose DevBench, a comprehensive benchmark that evaluates LLMs across various stages of the software development lifecycle, including software design, environment setup, implementation, acceptance testing, and unit testing. DevBench features a wide range of programming languages and domains, high-quality data collection, and carefully designed and verified metrics for each task. Empirical studies show that current LLMs, including GPT-4-Turbo, fail to solve the challenges presented within DevBench. Analyses reveal that models struggle with understanding the complex structures in the repository, managing the compilation process, and grasping advanced programming concepts. Our findings offer actionable insights for the future development of LLMs toward real-world programming applications. Our benchmark is available at https://github.com/open-compass/DevBench.

## 1. Introduction

Given its practical value and reasoning challenges, programming has become an important domain to deploy and evaluate large language models (LLMs), leading to popular products like GitHub Copilot and research benchmarks like HumanEval (Chen et al., 2021) and APPS (Hendrycks et al., 2021). While these earlier coding tasks focused on generating a single code file or even a single method conditioned on simple instructions, recent works such as SWE-bench (Jimenez et al., 2023) and RepoBench (Liu et al., 2023b) propose new evaluations that challenge LLMs with repository-level coding tasks, which feature longer, more involved NL2Code problems, that capture practical challenges in the real world. Still, these benchmarks concentrate on narrow aspects of software development. To this end, datasets and metrics that comprehensively evaluate software development as a multi-phase set of tasks are notably absent.

To address these shortcomings and fill this gap, we propose DevBench, a novel benchmark that mirrors real-world software development. DevBench generally evaluates models on the task of constructing a multi-file codebase starting from a product requirement document (PRD) of detailed specifications. Subscribing to the traditional Waterfall software development model (Royce, 1987), DevBench breaks down this process into a diverse set of inter-related development stages, i.e., software design, environment setup, implementation, acceptance and unit testing, as visualized in Figure 1 and Table 1. Across all tasks, DevBench is the first to evaluate models' software design and environment setup capabilities.

It is not trivial to collect a dataset for DevBench, as open-sourced repository data often lack essential components, especially the design documents and testing programs. The DevBench dataset is made up of 22 curated repositories in 4 languages (Python, C/C++, Java, JavaScript). These codebases cover a variety of domains, including machine learning, web service, command line utilities, etc. By encompassing the multi-faceted, interconnected steps of software development under a single framework, DevBench provides a holistic view of the capabilities of LLMs for automated software production, moving beyond the conventional, singular focus on code completion.

We also develop a baseline system by employing ChatDev (Qian et al., 2023a;b) as the foundational framework to benchmark state-of-the-art pretrained code models. We find that all models fail to solve the challenges presented

---

[*]Equal contribution  [1]Shanghai AI Laboratory [2]Nanjing University [3]Beijing University of Posts and Telecommunications [4]Dartmouth College [5]Princeton University [6]The University of Hong Kong [7]Tsinghua University [8]ByteDance Inc.. Correspondence to: Chao Peng <pengchao.x@bytedance.com>, Kai Chen <chenkai@pjlab.org.cn>. †: Work done during an internship at Shanghai AI Laboratory.
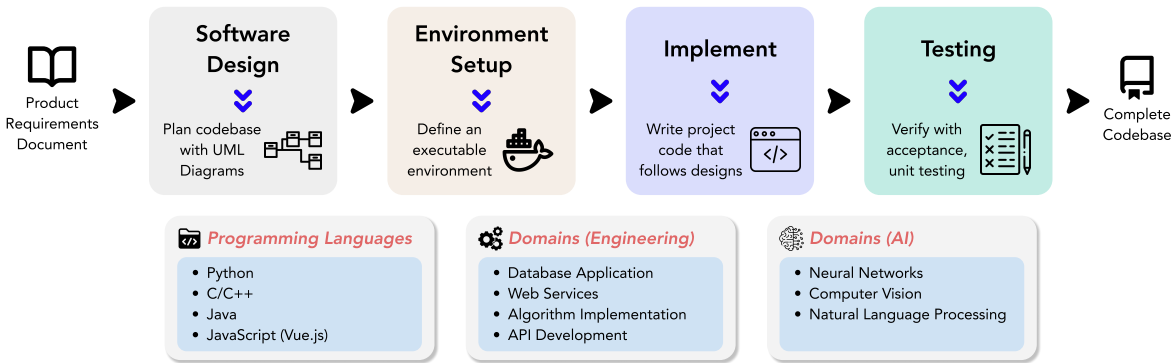
Preliminary work.

*Figure 1.* Our DevBench features multiple stages of software development, including software design, environment setup, implementation, and testing (both acceptance and unit testing).

| Task | Input† | Output | Environment | Evaluation |
|---|---|---|---|---|
| Software Design | PRD | UML Diagrams‡, Architecture Design | N/A | Subjective Evaluation |
| Environment Setup | PRD, UML Diagrams, Architecture Design | Dependency Files | Base | Pass Rate on Usage Examples |
| Implementation | PRD, UML Diagrams, Architecture Design | Implementation Code | Reference | Unit & Acceptance Testing |
| Acceptance Testing | PRD, UML Diagrams, Architecture Design, Implementation Code* | Acceptance Testing Code | Reference | Oracle Test |
| Unit Testing | PRD, UML Diagrams, Architecture Design, Implementation Code | Unit Testing Code | Reference | Oracle Test & Coverage |

*Table 1.* Task design in DevBench. †: following our modular evaluation protocol, the input for each task are reference. ‡: includes UML class and sequence diagrams. ∗: implementation source code is optional for acceptance testing.

within DevBench. GPT-4-Turbo achieves the highest scores amongst all evaluated models, yet it obtains less than 10% on our repository-level implementation task. Other tasks prove relatively more manageable for models; however, even GPT-4-Turbo struggles to attain scores exceeding 40% on any given task. For instance, models generally fail to generate executable testing code, with oracle test scores falling below 40%. Despite this, the generated testing code demonstrates potential in terms of coverage score, achiev- ing as high as 79.4% when it is executable. Furthermore, our investigation into different prompting methods shows that only prompts supplemented with external information, execution feedback in our context, yield notable and con- sistent improvements. We employ LLM-as-a-Judge (Zheng et al., 2023a) for evaluating software design, where the outcomes exhibit high agreement with human annotations. Importantly, qualitative analysis shows that models demon- strate difficulties in handling Makefile and Gradle, config- uring function arguments and employing advanced object- oriented programming techniques. Overall, DevBench in- troduces a novel challenge for existing code models, and

our investigation sheds light on fundamental issues, paving the way for future research.

## 2. Related Work

While initial code generation datasets mainly put forth self- contained, closed form completion problems predominantly written in Python (Chen et al., 2021; Hendrycks et al., 2021; Austin et al., 2021), subsequent works have raised coding tasks' complexity via a number of augmentations, including increasing language coverage (Zheng et al., 2023b; Cassano et al., 2023), expanding execution-based test coverage (Liu et al., 2023a; Wang et al., 2022), requiring use of dependen- cies (Lai et al., 2022; Ding et al., 2023; Liu et al., 2023c), proposing interactive coding environments (Yin et al., 2022; Yang et al., 2023), and constructing suites of short-form coding tasks (Lu et al., 2021; Muennighoff et al., 2023). De- vBench falls most in line with recent works that introduce *repository*-scale coding, where a task worker must gener- ate a code completion or patch that fixes an issue within a production-level codebase (Jimenez et al., 2023; Liu et al.,

2023b; Zhang et al., 2023). While these prior works mainly focus on iterative software development, they do not investigate LLMs' ability to create a codebase from extensive natural language project descriptions.

In automatic programming, recent works propose role-play frameworks for programming that involves communicative agents (Hong et al., 2023; Li et al., 2023; Qian et al., 2023a;b). Among them, MetaGPT (Hong et al., 2023) is the most closely related to our work, which integrates Standardized Operating Procedures into prompt sequences for streamlined workflows. Additionally, MetaGPT is characterized by structured output formats and a modular task design that resembles our work. The SoftwareDev auxiliary dataset they put forth is focused on cost, code statistics, and a general assessment of executability. In contrast, DevBench pioneers a novel evaluation approach founded on a carefully curated dataset of repositories. We define structured, thorough, and comprehensive evaluations across every phase of the software development lifecycle.

## 3. DevBench

In this section, we discuss the design of DevBench, including task specifications and the evaluation criteria (summarized in Table 1). Additional insights and exemplar artifacts of the tasks are presented in Appendix A. Adhering to the modular approach, our design utilizes reference inputs for each task, instead of relying on outputs generated by models from preceding stages. This strategy enables and concentrates on evaluating the efficacy of models in executing specific tasks. [1]

### 3.1. Task 1: Software Design

During this phase, the model is tasked with interpreting the Product Requirements Document (PRD) to create Unified Modeling Language (UML) diagrams, utilizing the extensively recognized Mermaid syntax, alongside formulating architecture designs using hierarchical file-tree structures. The class diagram acts as a fundamental component in software system modeling, illustrating classes, their attributes, and interrelations. Sequence diagrams detail processes and objects involved and the sequence of messages exchanged to carry out the software functionality. Architecture design is aimed at crafting a structured layout of source code, compilation scripts, and auxiliary files crucial for software implementation. For exemplifications of UML diagrams and architectural frameworks, refer to Appendix A.1.

---

[1] Notably, our framework can be adeptly configured to facilitate end-to-end evaluations, utilizing the intermediate outputs generated by models across multiple tasks. Moreover, DevBench can also function in a Copilot mode, empowering human users to intervene and refine model outputs, thus enhancing the collaborative synergy between human expertise and automated systems.

**Evaluation.** Since the Software Design tasks are open-ended, we employ the *LLM-as-a-judge* approach (Zheng et al., 2023a; Wang et al., 2023; Chiang & Lee, 2023) to conduct the automatic evaluation. The evaluation is anchored by two principal metrics: *general principles* and *faithfulness*. Detailed evaluation guidance is found in Appendix B.

The *general principles* metric plays a crucial role, with each task sharing common elements while maintaining specific criteria. For all the sub-tasks, principles like cohesion and decoupling, and practicability are fundamental. Cohesion and decoupling emphasize the importance of clarity and functionality within individual elements (classes or sequences) and reducing dependencies between different components. In terms of practicability, all tasks require designs to be readable, understandable, and modular, facilitating ease in development, testing, and maintenance. Meanwhile, each task has its unique focus areas: UML Class diagrams are evaluated on complexity; UML Sequence diagrams concentrate on uniformity, integration, and interaction complexity; Architecture Design highlights the distinction between design and coding, and conformance to practical standards. Subsequently, the *faithfulness* metric gauges the extent to which models adhere to specified instructions.

### 3.2. Task 2: Environment Setup

In the second phase of DevBench, models are provided with the PRD, UML diagrams, and architecture design to generate a dependency file requisite for initializing the development environment. This step is followed by the deployment of a standard installation command utilizing the generated file. For Python, the `Conda` environment manager is employed; for Java and JavaScript, `Gradle` and `NPM` are utilized respectively.[2] Setting up an environment often encounters challenges such as missing or outdated dependencies, along with version conflicts, all of which must be resolved to ensure a seamless development environment. Our research aims to investigate the potential of LLMs in automating this cumbersome process, thereby enhancing production efficiency.

**Evaluation.** The evaluation centers on the execution of dependency files across each programming language within a predetermined base environment delineated in a Docker file. This is followed by the execution of the repository's example usage code. The principal metric for evaluation in this task is the success rate of the executed example code.

---

[2] It is noteworthy that our evaluation framework does not contain an Environment Setup task for C/C++ due to the absence of a universally acknowledged and user-friendly dependency management system for these languages (Miranda & Pimentel, 2018).
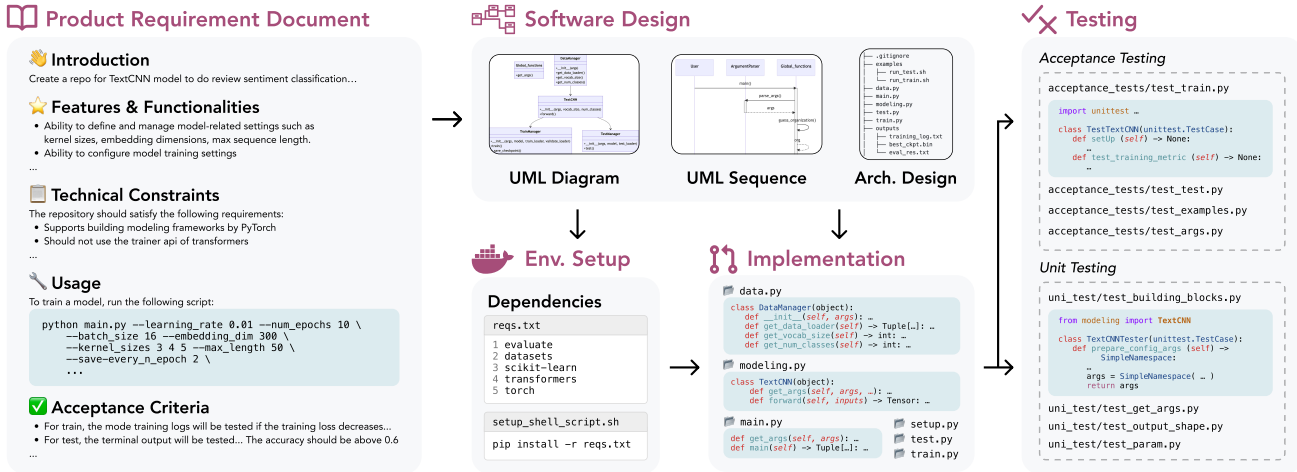
*Figure 2.* A typical workflow within DevBench showing building a codebase to develop a TextCNN model for movie review classification.

## 3.3. Task 3: Implementation

For this task, models are provided with the PRD, UML diagrams and architecture design, and are then instructed to develop code for each source code file as specified in the architecture design. Diverging from existing benchmarks on repository-level code generation (Liu et al., 2023b), the implementation task in DevBench is dedicated to assessing LLMs in generating an entire code repository from scratch. While similar in spirit to systems like MetaGPT (Hong et al., 2023) and ChatDev (Qian et al., 2023a;b), which derive output from brief requirement descriptions typically under 100 words, our approach is distinct. These systems indeed incorporate software design phases, yet the brevity and lack of detail in their input requirements pose significant challenges in evaluating the extent to which LLMs adhere to these initial specifications. In contrast, real-world software development scenarios often entail more detailed requirements to ensure the final product aligns precisely with expectations, with acceptance testing employed for verification (as discussed in the following subsection).

To more accurately simulate real-world development practices and ensure rigorous evaluation, the implementation task in DevBench involves supplying LLMs with comprehensive inputs including the PRD, UML class and sequence diagrams, and architecture design. The models are then prompted to generate code files. Given the constraint of output length, we adopt a sequential generation approach, prompting the models to produce one code file per interaction. Regarding the planning aspect, recent studies have explored prompting LLMs or training a specific planner (Yao et al., 2023; Besta et al., 2023; Wang et al.). Considering the structured nature of code files and the inherent dependencies among them, we utilize these dependencies as a clue for effective planning. The generation process is guided to adhere to a partial order derived from a predefined directed acyclic graph, thereby ensuring structured and logical code development. We leave the exploration of planning generated by models themselves for future work.

**Evaluation.** For the evaluation of the implementation task, an automated testing framework has been developed. This framework, tailored to the specific programming language in use, integrates `PyTest` for Python, `GTest` for C++, `JUnit` for Java, and `Jest` for JavaScript. The evaluation procedure involves executing reference acceptance and unit tests within a predefined reference environment and the evaluation metric is determined by the pass rate of these tests.

## 3.4. Task 4: Acceptance Testing

For this task, models are provided with the PRD, UML diagrams, and architecture design, with the option to include the implementation source code to generate acceptance test code. Acceptance testing is critical to verify that the software adheres to requirements and operates effectively. In the context of applications featuring command-line interfaces, acceptance tests interact with the software via shell commands, as specified in the PRD, and subsequently evaluate the accuracy of the output generated. For libraries, acceptance tests are implemented through code that invokes the library's API, followed by assertions made on the responses of the API. Applying this evaluative approach to LLMs provides valuable insights into their practical effectiveness and dependability within the domain of software development.

**Evaluation.** The evaluation of this task involves running the generated acceptance tests against the benchmark implementation code in the same testing framework developed for the evaluation for the implementation task.

| | **Python** | **C/C++** | **Java** | **JavaScript**[‡] |
|---|---|---|---|---|
| Domain[†] | NLP  CV  DL  ALGO  API  Tool | DB  ALGO  Tool | CV  DB  ALGO  Tool | WEB |
| #Repo | 10 | 5 | 5 | 2 |
| Avg. #Code File | 2.2 | 7.0 | 5.4 | 6.0 |
| Avg. #Code Line | 276 | 495 | 524 | 617 |
| Avg. #Accep. Tests | 3.0 | 5.4 | 2.4 | 2 |
| Avg. #Unit Tests | 12.4 | 11.8 | 8.2 | - |
| Avg. Coverage | 91.8 | 95.0 | 64.9* | - |

*Table 2.* DevBench Statistics. †: DevBench covers a range of domains including NLP, computer vision, deep learning, algorithm implementation, API applications, Database applications, web service (both frontend and backend), and general tools and utilities. ‡: We do not include unit testing for JavaScript as checking functional correctness is not applicable to pure static web pages. Correctness of page rendering and user interaction handling is checked using acceptance tests. ∗: Interfaces with thirty-party libraries that are not used to implement the designated functionalities are not supplied with test cases, resulting in relatively lower overall test coverage.

In this phase, the Oracle Test methodology is employed to assess the accuracy of acceptance tests generated by the model, which are comprised of both test input and expected output. The testing framework conducts an execution of the reference implementation of the software using the input devised by the model, subsequently comparing the software's actual output against the model-predicted output. This approach facilitates an understanding of the extent to which LLMs can accurately interpret and predict the correct behavior of the subject software, as delineated in its design documentation.

### 3.5. Task 5: Unit Testing

In this phase, the PRD, UML diagrams, and architecture design are furnished to the models to generate unit tests, facilitating a comprehensive understanding of the software. Unit testing serves as a fundamental approach to safeguard code integrity and operational accuracy. Distinguished from broader acceptance testing, unit testing examines individual code segments for adherence to specified functionalities. The increasing reliance on LLMs for streamlining software development processes underscores the criticality of their adeptness in writing effective unit tests, a competency integral to the software's reliability and overall robustness.

**Evaluation.** The evaluation of LLM-generated unit tests is

conducted through their application on the reference source code, employing the previously delineated testing framework. Similar to the acceptance testing evaluation, this phase leverages the Oracle Test, wherein the actual output of code units under specific test inputs is compared to anticipated outputs, as delineated by the oracle.

In addition, code coverage metrics are incorporated, providing a quantitative understanding of test comprehensiveness. Utilizing statement coverage analysis tools integrated within the aforementioned testing frameworks, coverage is mathematically expressed as:

$$\text{Coverage} = \left( \frac{\text{Number of Executed Statements}}{\text{Total Number of Statements}} \right) \times 100\%,$$

where the number of executed statements denotes the count of distinct executable statements within the code that are executed at least once during the testing process, while the total number of statements represents the aggregate count of all executable statements present in the codebase that are subject to potential execution.

| Task | Implementation | |
|---|---|---|
| Evaluation Metric (%) | Pass@ Accept. Test[¶] | Pass@ Unit Test[¶] |
| GPT-4-Turbo | | |
| No-Review | 3.0 | 0.0 |
| Normal-Review | 3.0 | 0.0 |
| Execution-Feedback | 8.9 | 4.2 |
| CodeLlama-34B-Instruct | | |
| No-Review | 0.0 | 1.4 |
| Normal-Review | 0.0 | 1.4 |
| Execution-Feedback | 0.0 | 1.4 |
| DeepSeek-Coder-33B-Instruct | | |
| No-Review | 1.5 | 4.2 |
| Normal-Review | 1.5 | 4.2 |
| Execution-Feedback | 1.5 | 4.2 |

*Table 3.* Results of different prompting methods of the Implementation task on a subset of DevBench. ¶: all results are averaged across all repositories and weighted by the number of code lines, which measures the difficulty of each repository.

## 4. Dataset

### 4.1. Dataset Construction

The data preparation process in our work consists of three distinct phases: repository preparation, code cleanup, and document preparation.

The initial phase, repository preparation, involves selecting high-quality, well-structured candidate repositories from a GitHub dump. Recognizing the impracticality of constructing a repository from scratch, we employed a filtering process to identify suitable candidates. Moreover, we im-

posed a constraint on the total number of lines of code to ensure the repositories' complexity remains manageable, facilitating the evaluation of current LLMs.

During the code cleanup phase, our postgraduate student annotators were tasked with setting up the required environment as stipulated in the repositories' README files. They then executed the code to verify its functionality. Following this sanity check, the annotators were instructed to meticulously refine the code repositories. This refinement included the removal of unnecessary auxiliary files. To ascertain code quality, the annotators were also required to run existing unit and acceptance tests, or to develop additional tests, ensuring they meet the standards of the oracle test and achieve satisfactory coverage.

The final phase, document preparation, involved the creation of standard software design documents, namely, UML Class and Sequence Diagrams, and Architecture Designs, for each repository. We provide annotators with specific guidelines and templates for these documents. The annotators were responsible for ensuring that these design documents corresponds accurately and cohesively with the respective code repositories.

### 4.2. Dataset Statistics

In Table 2, we present an exhaustive statistical breakdown of our datasets. DevBench contains a collection of 22 curated repositories, spanning across four widely-used programming languages (Python, C/C++, Java, JavaScript) and a diverse range of domains. For detailed repositories statistics, refer to Table 8 in Appendix D. The dataset is characterized by its multi-file structure. The Python repositories in our dataset are relatively straightforward, with each repository comprising approximately two files and an average of 276 lines of code. In contrast, the repositories pertaining to statically-typed programming languages, namely C/C++ and Java, are more complex, featuring an increased count of code files and lines. For JavaScript, our usage of the Vue.js framework necessitates that models adeptly navigate the framework's templates and development paradigms. Consequently, JavaScript repositories exhibit the highest number of code files and lines, posing a substantial challenge for LLMs. Additionally, we have prepared extensive reference acceptance and unit tests for each repository to facilitate rigorous evaluation of the implementation task.

## 5. A Baseline System

We introduce a baseline system formulated for our DevBench, building upon the foundations of ChatDev (Qian et al., 2023a;b). ChatDev is a virtual, chat-powered software development system that adheres to the conventional waterfall model. It bifurcates the development process into

four primary tasks (dubbed as phrases in ChatDev): design, coding, testing, and documentation. Within this system, multiple LLM agents assume diverse roles such as programmers, reviewers, and testers, pertinent to each phase. ChatDev is characterized by its utilization of a chat chain mechanism, which segments each phase into smaller, atomic tasks. This approach enables context-sensitive, multi-turn dialogues between two distinct roles, facilitating the proposal and validation of solutions for individual tasks.

In contrast to ChatDev, our development of baseline system incorporates several features and enhancements. We have restructured the task design to align closely with the evaluation criteria of DevBench. Specifically, this includes the integration of comprehensive input to the system, exemplified by well-structured PRDs. This integration is crucial in addressing and examining the issue of hallucination in controlled experimental settings. Moreover, our baseline system expands upon the capabilities of ChatDev, supporting a wider range of tasks, including standard Object-Oriented programming designs (UML Class and Sequence diagrams), repository planning (architecture design), environment setup, and acceptance testing. A significant advancement in our baseline system is its compatibility with multiple programming languages and their corresponding runtime environments. This feature is coupled with the provision of comprehensive execution feedback to the system.

| Task | Accept. Testing |
|---|---|
| Evaluation Metric (%) | Oracle Test[§] |
| GPT-4-Turbo | |
|   No-Review | 3.6 |
|   Normal-Review | 7.7 |
|   Normal-Review w/ Src Code | 14.9 |
| CodeLlama-34B-Instruct | |
|   No-Review | 0.0 |
|   Normal-Review | 0.0 |
|   Normal-Review w/ Src Code | 0.0 |
| DeepSeek-Coder-33B-Instruct | |
|   No-Review | 0.0 |
|   Normal-Review | 4.2 |
|   Normal-Review w/ Src Code | 15.6 |

*Table 4.* Results of different prompting methods of the Acceptance Testing task on a subset of DevBench. §: all results are averaged across all repositories and weighted uniformly.

## 6. Experiments

### 6.1. Setup

**Base Models** We evaluate three prominent pre-trained model families with different model sizes, including both proprietary and open-source models: OpenAI GPT (GPT, 2023), CodeLlama (Roziere et al., 2023), DeepSeek-

Coder (Guo et al., 2024). Specifically, our experiments involve GPT-3.5-Turbo, GPT-4-Turbo from OpenAI GPT[3], CodeLlama-Instruct 7B/13B/34B for CodeLlama, and DeepSeek-Coder-Instruct models 1.3B/6.7B/33B for DeepSeek-Coder. For comprehensive details on these models, we refer readers to the original papers.

**Prompting Methods** In our baseline system, we explore three prompting methods: `No-Review`, `Normal-Review`, and `Execution-Feedback`. `No-Review` represents a basic zero-shot prompting with built-in task prompts. `Normal-Review` involves a dual-role interaction, where the first role generates a solution and the second role reviews and, where necessary, corrects it. This mode is designed to evaluate the impact of review on model performance, in the absence of external inputs. `Execution-Feedback`, on the other hand, adds more dynamic interaction to the review process. This feedback includes runtime results, error messages, and performance metrics. Such information could enable the reviewing role to make more informed decisions, potentially leading to more accurate and effective solutions. To save API and computational costs, we only review once for all review-involved prompting methods.

**Implementation Details** We utilize LMDeploy for the deployment of CodeLlama and DeepSeek-Coder models.[4] Acknowledging the potential for extensive input context in DevBench tasks, we configure the context length to 32K for these models. For the Software Design task, we set the temperature parameter to 0.2, while for the remaining four tasks, we use a temperature of 0. Other hyperparameters in the experiment are maintained at default settings. All code-related tasks are rigorously evaluated in an isolated sandbox environment, utilizing Docker technology.

**Evaluation on Software Design** We follow previous work (Zheng et al., 2023a) to conduct a pairwise comparison to determine which response is better, focusing on the metrics of *general principles* and *faithfulness* (see the details in Appendix B). To reduce the expenditure of the OpenAI GPT API and human effort, the scope of our evaluation was confined to a subset of our Benchmark. This process involves 192 pairs across eight repositories, eight models, and three sub-tasks. Regarding the LLM judge, we use GPT-4-Turbo as the judge and GPT-3.5-Turbo as the baseline model. To mitigate the issue of position bias (i.e., preferring response at a certain position regardless of the content), the evaluation was executed in a *dual* mode, evaluating each pair twice in different orders (384 pairs in total), with inconsistent decisions being considered as a tie. For the human evaluation, we shuffle the order of two responses

---

[3] We utilize `gpt-3.5-turbo-1106`, `gpt-4-0125-preview`, respectively.

[4] https://github.com/InternLM/lmdeploy

and annotate each pair thrice to obtain the *human majority*.

| Task | Unit Testing | |
|---|---|---|
| Evaluation Metric (%) | Oracle Test[§] | Coverage[$] |
| GPT-4-Turbo | | |
| No-Review w/ Src Code | 35.1 | 34.3 (54.8) |
| Normal-Review w/ Src Code | 22.6 | 25.8 (68.7) |
| CodeLlama-34B-Instruct | | |
| No-Review w/ Src Code | 10.7 | 18.3 (73.2) |
| Normal-Review w/ Src Code | 12.6 | 27.0 (72.1) |
| DeepSeek-Coder-33B-Instruct | | |
| No-Review w/ Src Code | 27.2 | 37.9 (75.8) |
| Normal-Review w/ Src Code | 22.5 | 35.5 (71.0) |

*Table 5.* Results of different prompting methods of the Unit Testing task on a subset of DevBench. §: the Oracle Test results are averaged across all repositories and weighted uniformly. $: the results on the left side are averaged across all repositories and weighted uniformly, showing the overall scores. The results on the right side in the parenthesis are averaged across all *valid* repositories and weighted uniformly, where models have generated *executable* testing code.

## 6.2. Main Results

**Results across prompting methods** We first conduct experiments to examine the effects of different prompting methods on a subset of DevBench using three representative models: GPT-4-Turbo, CodeLlama-34B-Instruct, and DeepSeek-Coder-33B-Instruct.

Table 3 illustrates the results of the implementation task for our study. In general, `Execution-Feedback` leads to the optimal performance, where GPT-4-Turbo benefits the most, especially on acceptance tests. In contrast, CodeLlama-34B-Instruct and DeepSeek-Coder-33B exhibited no improvement with any review process. We note that the efficacy of the review process could be understated as our automated testing is too rigorous and sparse to reflect the improvements. We observe that, despite no substantial improvements on reference tests, the code quality notably improved with `Execution-Feedback` prompt. However, there is no significant improvements using the `Normal-Review` setting compared with the `No-Review` setting. Models consistently provide encouraging feedback on generated code and offer unhelpful suggestions, such as the addition of unnecessary error handling or reorganization. This indicates that the models are unable to comprehend complicated code by merely reading it, lacking external knowledge or assistance like execution feedback.

For the testing tasks, which are relatively easier than the implementation, there are various observations. In Table 5, we find no clear evidence that the review process brings stable benefits to unit testing. However, on the acceptance testing,

| Task | Environment Setup | Implementation | | Acceptance Testing | Unit Testing | |
|---|---|---|---|---|---|---|
| Evaluation Metric (%) | Pass@ Example Usage$^\S$ | Pass@ Accept. Test$^\P$ | Pass@ Unit Test$^\P$ | Oracle Test$^\S$ | Oracle Test$^\S$ | Coverage$^\$$ |
| GPT-3.5-Turbo | *33.3* | 4.2 | 4.3 | 11.7 | 28.7 | 24.6 (61.4) |
| GPT-4-Turbo | ***41.7*** | **7.1** | **8.0** | **29.2** | **36.5** | 33.2 (66.3) |
| CodeLlama-7B-Instruct | *8.3* | 0.0 | 0.0 | 0.0 | 3.0 | 3.6 (71.0) |
| CodeLlama-13B-Instruct | *25.0* | 0.6 | 0.0 | 0.0 | 5.1 | 8.6 (57.6) |
| CodeLlama-34B-Instruct | *16.7* | 0.6 | 0.5 | 4.5 | 21.1 | 25.4 (72.6) |
| DeepSeek-Coder-1.3B-Instruct | *8.3* | 0.0 | 0.1 | 0.0 | 5.6 | 2.7 (27.0) |
| DeepSeek-Coder-6.7B-Instruct | *25.0* | 2.9 | 3.9 | 20.5$^\heartsuit$ | 23.5 | 28.2 (70.6) |
| DeepSeek-Coder-33B-Instruct | *16.7* | 4.4 | 5.5 | 13.6 | 32.8 | 35.7 (79.4) |

*Table 6.* Task 2 to Task 5 results on DevBench. *Italic figures*: test cases for the Environment Setup task are quite scarce compared to other tasks, therefore the results are more influenced by the randomness.[6] $\S$: all results are averaged across all repositories and weighted uniformly. $\P$: all results are averaged across all repositories and weighted by the number of code lines. $\$$: the results on the left side are averaged across all repositories and weighted uniformly, showing the overall scores. The results on the right side in the parenthesis are averaged across all *valid* repositories and weighted uniformly, where models have generated *executable* testing code. $\heartsuit$: the model has generated meaningless but executable testing code.

`Normal-Review` brings enhancement to the performance. Regarding the visibility of implementation source code, it is common practice not to expose source code and execution feedback for acceptance testing, while it's allowed to employ the source code as input for unit testing. As shown in Table 4, models can barely generate executable acceptance testing code and incorporating source code as additional input dramatically increases the performance.

**Results across base models** Table 6 illustrates the main results on DevBench with optimal prompting methods applied for each task.[5] We find that GPT-4-Turbo demonstrates superior performance compared to other models, while all models evaluated are far from satisfactory. DevBench can effectively distinguish between models of varying capabilities. Smaller models, such as CodeLlama-7B/13B-Instruct and DeepSeek-Coder-1.3B-Instruct, demonstrate inherent limitations, frequently unable to generate syntactically accurate code or follow the instructions. These models tend to generate mere code skeletons or fill the function body with only comments. Larger open-sourced models and GPT models, while generating more reasonable code, still struggle with the subtleties of complex code structure and logic, such as variable type conversion, function arguments and object-oriented class.

**Results across tasks** Generally, the models' performances on the implementation task are all below the 10% pass rate. The highest-performing model, GPT-4-Turbo, registers only a 7.1% pass rate on reference acceptance tests and 8.0% on unit tests, while some other models score zero, failing all reference tests. Despite prior research such as HumanEval (Chen et al., 2021) indicating that models can manage simple code-writing tasks, substantial challenges remain in more complex coding scenarios within DevBench. We break down the implementation results into different languages in Figure C and find that models particularly struggle in handling Java and C/C++, whose stringent syntax requirements tend to magnify the models' deficiencies in managing intricate details. This points to the necessity for more enriched and diverse training data across programming languages to bridge this competency gap.

Compared with the implementation task, other tasks are relatively simple but still challenging. Regarding the environment setup task, we note that the test cases for the environment setup task are quite scarce compared to other tasks[6], therefore the results are more influenced by the randomness. Roughly speaking, GPT-4-Turbo reaches a 41.7% pass rate in building environments, while open-sourced models largely fall behind it. For the testing tasks, GPT-4-Turbo still obtains the highest scores and open-sourced models perform worse. We identify an outlier that DeepSeek-Coder-6.7B-Instruct achieves the highest score among open-sourced models and approach GPT-4-Turbo in acceptance testing. The model somehow *cheats* on this task by generating meaningless but executable testing code (see Appendix C.4). With respect to the unit testing, the overall Oracle Test and Coverage scores are quite low, which are averaged across all repositories. However, for those generated testing codes that

---

[5] `Execution-Feedback` is used for Environment Setup and Implementation; `Normal-Review w/ Src Code` for Acceptance Testing; `No-Review w/ Src Code` for Unit Testing as `Normal-Review w/ Src Code` shows no clear advantage.

[6] Except for C/C++ and easy Python repositories that are absent, only 12 repositories are involved in the environment setup task. We will resolve this issue in future work. However, DevBench contains rich test cases for other tasks. Table 2 evidences that DevBench features fruitful tests for the implementation task. For the testing generation tasks, our prompts depict fine-grained requirements and ensure the quantity of generated testing cases.

can be successfully executed (obtain non-zero Oracle Test score), the coverage scores are relatively high, suggesting the models' promising capability on this problem.

**Results on software design** Table 7 shows the results of software design using GPT-4-Turbo as the Judge and GPT-3.5-Turbo as the baseline reference. GPT-4-Turbo dominantly outperforms GPT-3.5-Turbo with extraordinarily high win rates in all cases. Regarding the open-sourced models, as the size increases, models consistently produce higher quality design documents on both metrics, while they are relatively inferior on *faithfulness*.

We compare the LLM-as-a-Judge results with human majority annotations. Low agreements are observed with tie considered, which aligns with the findings in previous studies (Zheng et al., 2023a). It is reasonable as a tie is hard to define and judge, especially for highly complicated and structured software design documents. Without tie, GPT-4-Turbo reaches 79.2% and 83.2% agreements on the *general principles* and *faithfulness* metrics, respectively. This means GPT-4-Turbo's judgments align with the majority of humans and could serve as a good alternative for automated software design evaluation.

| | w/ Tie | | w/o Tie | |
|---|---|---|---|---|
| | G† | F‡ | G | F |
| GPT-4-Turbo | **97.9** | **97.9** | **100.0** | **100.0** |
| CodeLlama-7B-Instruct | 4.2 | 8.3 | 4.2 | 4.5 |
| CodeLlama-13B-Instruct | 18.8 | 14.6 | 10.5 | 5.3 |
| CodeLlama-34B-Instruct | 39.6 | 33.3 | 33.3 | 21.4 |
| DeepSeek-Coder-1.3B-Instruct | 16.7 | 16.7 | 5.5 | 5.6 |
| DeepSeek-Coder-6.7B-Instruct | 35.4 | 35.4 | 31.6 | 29.4 |
| DeepSeek-Coder-33B-Instruct | 52.1 | 50.0 | 53.8 | 50.0 |
| Agree w/ Human Majority | 60.4 | 51.6 | 79.2 | 83.2 |

*Table 7.* Win rate of pairwise comparison against GPT-3.5-Turbo on Software Desgin on a subset of DevBench where results are averaged across different repositories and sub-tasks uniformly. †: the *general principles* metric. ‡: the *faithfulness* metric. w/ Tie: inconsistent results are considered as a tie. We also report agreement with Human Majority.

### 6.3. Analysis

**Models struggle with Makefile and Gradle.**. LLMs often face challenges in generating accurate Makefile for C/C++ and Gradle files for Java. Frequently, the generated files are deficient in critical components like source code files, necessary dependencies and essential tasks. In C/C++ and Java repositories, approximately 90% of compilation and execution errors can be attributed to these issues. We find that even GPT-4-Turbo occasionally fail in basic syntax errors on compilation files. This is potentially caused by insufficient training data related to these compilation tools.

**Models fail to configure function arguments**. We observe that models have difficulties in accurately following complex requirements, often failing to correctly configure and utilize function arguments. This contrasts with previous studies that highlight the models' relative proficiency in implementing function bodies. Such errors in arguments are detected in 20% of the repositories examined using GPT-4-Turbo, in 30% with CodeLlama-34B-Instruct and in 20% with DeepSeek-Coder-33B-Instruct. These issues cannot be effectively corrected through our review process with either `Normal-Review` or `Execution-Feedback` settings. This is because models frequently misattribute execution errors to incorrect usage of test programs, despite the clear prompt that the gold testing has already been given.

**Models face obstacles in advanced programming skills.**. Another observation reveals that models encounter difficulties with advanced programming skills such as overloading, utilizing pointers, and object-oriented programming. In 20% of the repositories tested, there is a noticeable misunderstanding of such advanced notions. Models are unable to distinguish between private and public members/methods, often erroneously attempting to access private variables or import class functions into other code files. Despite receiving detailed feedback on execution errors, models still fail to correctly address these issues, indicating a knowledge gap in complex programming features and skills.

## 7. Conclusion

The DevBench framework introduced in this paper presents a leap forward in evaluating LLMs within the domain of automated software development. By employing a multi-stage benchmarking process, DevBench comprehensively assesses LLMs across a spectrum of tasks including design, environment setup, implementation, and testing. Empirical findings reveal that current pre-trained models like GPT-4-Turbo are still confronted with substantial challenges within DevBench. Through analysis, we identify models' limitations in understanding the complex repository structures and handling the nuanced demands of comprehensive software development. These insights obtained from the DevBench assessments elucidate critical pathways for future code model development.

## Impact Statement

DevBench and prior works demonstrate the exciting potential and impact that LLM oriented methods can have towards automating and assisting software engineering. As more automatic prototypes and processes will likely emerge to address different parts of various software engineering workflows, it becomes increasingly important to understand models' performance on not only each individual task, but

also how model output on one task affects subsequent downstream tasks that then receive the generation as input. Especially as LLM powered frameworks are deployed in real-world scenarios, robust and realistic model evaluation is critical to answering whether their designs and implementations are ready to be used and affect human end users. Our work's novel, holistic evaluation of software engineering incorporates tasks is an important first step towards answering such questions, and we are hopeful that DevBench may inspire more work that covers additional areas of software engineering where LLMs could play a role.

In addition to the task formulations and experimental findings presented in the main paper, we also provide detailed descriptions, motivations, examples, data, and additional results in the Appendix to supplement our claims and demonstrate how other practitioners can easily create their own code repositories and system design documents for evaluation using DevBench. While DevBench's current 22 repositories provide strong coverage of popular programming languages and use cases, we are eager to continually expand this coverage over time.

# References

Openai gpt, 2023. URL https://platform.openai.com/docs/models/overview.

Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., and Sutton, C. Program synthesis with large language models, 2021.

Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Gianinazzi, L., Gajda, J., Lehmann, T., Podstawski, M., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.

Cassano, F., Gouwar, J., Nguyen, D., Nguyen, S. D., Phipps-Costin, L., Pinckney, D., Yee, M.-H., Zi, Y., Anderson, C. J., Feldman, M. Q., Guha, A., Greenberg, M., and Jangda, A. Multipl-e: A scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions on Software Engineering*, 49:3675–3691, 2023. URL https://api.semanticscholar.org/CorpusID:258205341.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

Chiang, C.-H. and Lee, H.-y. Can large language models be an alternative to human evaluations? *ArXiv preprint arXiv:2305.01937*, 2023.

Ding, Y., Wang, Z., Ahmad, W. U., Ding, H., Tan, M., Jain, N., Ramanathan, M. K., Nallapati, R., Bhatia, P., Roth, D., and Xiang, B. Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion. *ArXiv*, abs/2310.11248, 2023. URL https://api.semanticscholar.org/CorpusID:264172238.

Guo, D., Zhu, Q., Yang, D., Xie, Z., Dong, K., Zhang, W., Chen, G., Bi, X., Wu, Y., Li, Y. K., Luo, F., Xiong, Y., and Liang, W. Deepseek-coder: When the large language model meets programming – the rise of code intelligence, 2024.

Hendrycks, D., Basart, S., Kadavath, S., Mazeika, M., Arora, A., Guo, E., Burns, C., Puranik, S., He, H., Song, D., and Steinhardt, J. Measuring coding challenge competence with apps. *NeurIPS*, 2021.

Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Zhang, C., Wang, J., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., and Schmidhuber, J. Metagpt: Meta programming for a multi-agent collaborative framework, 2023.

Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.

Lai, Y., Li, C., Wang, Y., Zhang, T., Zhong, R., Zettlemoyer, L., Yih, S., Fried, D., yi Wang, S., and Yu, T. Ds-1000: A natural and reliable benchmark for data science code generation. *ArXiv*, abs/2211.11501, 2022. URL https://api.semanticscholar.org/CorpusID:253734939.

Li, G., Hammoud, H. A. A. K., Itani, H., Khizbullin, D., and Ghanem, B. Camel: Communicative agents for" mind" exploration of large scale language model society. *arXiv preprint arXiv:2303.17760*, 2023.

Liu, J., Xia, C. S., Wang, Y., and Zhang, L. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *arXiv preprint arXiv:2305.01210*, 2023a.

Liu, T., Xu, C., and McAuley, J. Repobench: Benchmarking repository-level code auto-completion systems. *arXiv preprint arXiv:2306.03091*, 2023b.

Liu, Y., Tang, X., Cai, Z., Lu, J., Zhang, Y., Shao, Y., Deng, Z., Hu, H., Yang, Z., An, K., Huang, R., Si, S., Chen, S., Zhao, H., Li, Z., Chen, L., Zong, Y., Wang, Y., Liu, T., Jiang, Z., Chang, B., Qin, Y., Zhou, W., Zhao, Y., Cohan, A., and Gerstein, M. B. Ml-bench: Large language models leverage open-source libraries

for machine learning tasks. *ArXiv*, abs/2311.09835, 2023c. URL https://api.semanticscholar.org/CorpusID:265221105.

Lu, S., Guo, D., Ren, S., Huang, J., Svyatkovskiy, A., Blanco, A., Clement, C., Drain, D., Jiang, D., Tang, D., et al. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.

Miranda, A. and Pimentel, J. On the use of package managers by the c++ open-source community. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing*, pp. 1483–1491, 2018.

Muennighoff, N., Liu, Q., Liu, Q., Zebaze, A., Zheng, Q., Hui, B., Zhuo, T. Y., Singh, S., Tang, X., von Werra, L., and Longpre, S. Octopack: Instruction tuning code large language models. *ArXiv*, abs/2308.07124, 2023. URL https://api.semanticscholar.org/CorpusID:260886874.

Qian, C., Cong, X., Liu, W., Yang, C., Chen, W., Su, Y., Dang, Y., Li, J., Xu, J., Li, D., Liu, Z., and Sun, M. Communicative agents for software development, 2023a.

Qian, C., Dang, Y., Li, J., Liu, W., Chen, W., Yang, C., Liu, Z., and Sun, M. Experiential co-learning of software-developing agents. *arXiv preprint arXiv:2312.17025*, 2023b.

Royce, W. W. Managing the development of large software systems: concepts and techniques. In *Proceedings of the 9th international conference on Software Engineering*, pp. 328–338, 1987.

Roziere, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X. E., Adi, Y., Liu, J., Remez, T., Rapin, J., et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

Wang, P., Li, L., Chen, L., Zhu, D., Lin, B., Cao, Y., Liu, Q., Liu, T., and Sui, Z. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*, 2023.

Wang, Z., Cai, S., Liu, A., Ma, X., and Liang, Y. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*.

Wang, Z., Zhou, S., Fried, D., and Neubig, G. Execution-based evaluation for open-domain code generation. In *Conference on Empirical Methods in Natural Language Processing*, 2022. URL https://api.semanticscholar.org/CorpusID:254877069.

Yang, J., Prabhakar, A., Narasimhan, K., and Yao, S. Intercode: Standardizing and benchmarking interactive coding with execution feedback. *ArXiv*, abs/2306.14898, 2023. URL https://api.semanticscholar.org/CorpusID:259262186.

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023.

Yin, P., Li, W.-D., Xiao, K., Rao, A. E., Wen, Y., Shi, K., Howland, J., Bailey, P., Catasta, M., Michalewski, H., Polozov, O., and Sutton, C. Natural language to code generation in interactive data science notebooks. *ArXiv*, abs/2212.09248, 2022. URL https://api.semanticscholar.org/CorpusID:254854112.

Zhang, F., Chen, B., Zhang, Y., Liu, J., Zan, D., Mao, Y., Lou, J.-G., and Chen, W. Repocoder: Repository-level code completion through iterative retrieval and generation. In *Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://api.semanticscholar.org/CorpusID:257663528.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023a.

Zheng, Q., Xia, X., Zou, X., Dong, Y., Wang, S., Xue, Y., Wang, Z., Shen, L., Wang, A., Li, Y., et al. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*, 2023b.

# A. Software Engineering Tasks with a Running Example

We describe concepts of software engineering tasks using one of the subjects of DevBench, named Actor Relationship Game. The Actor Relationship Game is a Java-based application that allows users to explore connections between popular actors through their movie collaborations, using data from The Movie Database (TMDB) API. It constructs an actor graph and identifies the shortest path of relationships between any two actors.

## A.1. Software Design

Software design is the process by which an agent creates a specification of a software artifact, intended to accomplish goals, using a set of primitive components and subject to constraints. It is a phase in the software development lifecycle that bridges the gap between software requirements analysis and the actual implementation of the software system.

During the software design phase, software engineers or designers define the way a software application will work to meet the specified requirements in the form of a Product Requirement Document (PRD). They create diagrams that determine the data structures, software architecture, interface designs, and module specifications with Unified Markup Language (UML) diagrams.

A good software design is crucial as it impacts the quality, maintainability, performance, scalability, and robustness of the software product. It facilitates a smoother implementation phase, allows for better understanding and communication among team members, and helps in identifying potential issues early in the development process.

### A.1.1. CLASS DIAGRAMS

Class diagrams are a cornerstone of object-oriented design, offering a static snapshot of the system structure. These diagrams illustrate the classes within the system, their attributes, methods, and the relationships among the classes, such as inheritance and associations. Class diagrams are instrumental in providing an abstract representation of the system's components and their interactions, facilitating a deeper understanding of the software's overall architecture and design patterns.

Figure A.1.1 shows the class diagram of the Actor Relationship Game repository. This UML class diagram delineates the architecture of a system designed to model and analyze the network of relationships between actors and movies through an `Actor Graph`. It encompasses classes such as 'Actor' and 'Movie' to represent individual entities, alongside an `ActorGraph` class that serves as a repository and management layer for these entities and their associations. Utility and operational classes like `ActorGraphUtil`, `GameplayInterface`, and `GraphCreation` provide mechanisms for manipulating the graph—ranging from data ingestion using the `TMDBApi` to utility functions and gameplay interfaces that leverage the graph for various applications. The relationships between classes, including associations, aggregations, and dependencies, are meticulously outlined to depict interactions such as actors appearing in movies and the construction and utilization of the actor-movie graph for finding connections and supporting gameplay or analysis tasks.

### A.1.2. SEQUENCE DIAGRAMS

Sequence diagrams, different than class diagrams, focus on the dynamic aspects of the system. They depict how objects interact with each other across time, outlining the sequence of messages exchanged between objects to accomplish a specific functionality or process within the system. Sequence diagrams are invaluable for visualizing and analyzing the flow of operations, timing constraints, and the interaction patterns among system components, making them essential for detailed behavioral analysis.

The sequence diagram of `Actor Relationship Game` repository is illustrated in Figure A.1.2. This diagram illustrates the flow of operations for creating, populating, and utilizing an actor-movie graph. Initially, the `Main` function triggers the graph creation process by calling `createGraph()` on the `GraphCreation` module, which then interacts with the `TMDBApi` to fetch popular actors' data. Upon receiving this data, `GraphCreation` populates the `ActorGraph` with actors, movies, and their associations. After constructing the graph, `GraphCreation` delegates the responsibility of saving this graph to a file to `ActorGraphUtil`, which then returns a serialized file (`actorGraph.ser`) back to `Main`. Subsequently, `Main` instructs the `GameplayInterface` to load this graph and use it to find connections between actors via `findConnectionWithPath()`, a method in `ActorGraph`. The path found is then returned to `GameplayInterface`, which finally displays the results back in the `Main` function. This sequence encapsulates a complete lifecycle from graph creation, through data population and serialization, to utilization for finding actor connections, showcasing a systematic approach to managing and analyzing actor-movie relationships.
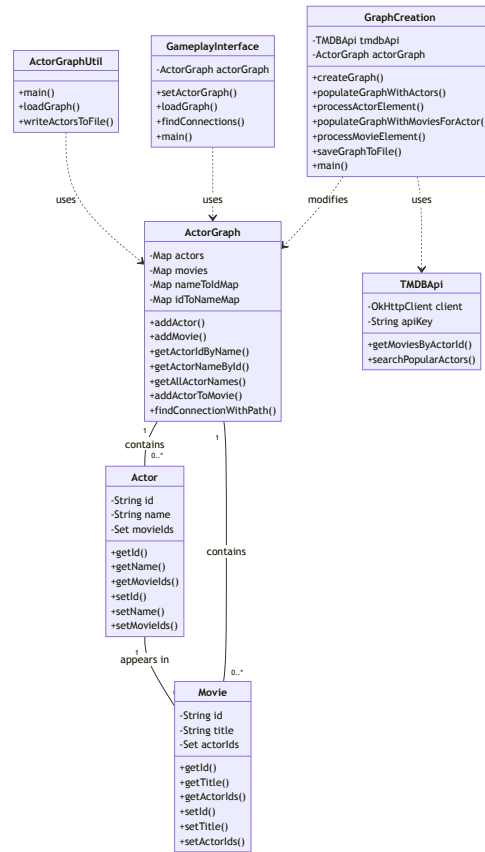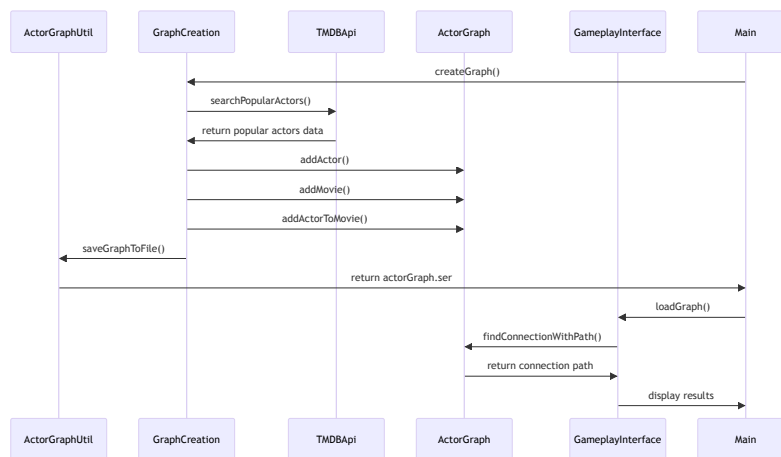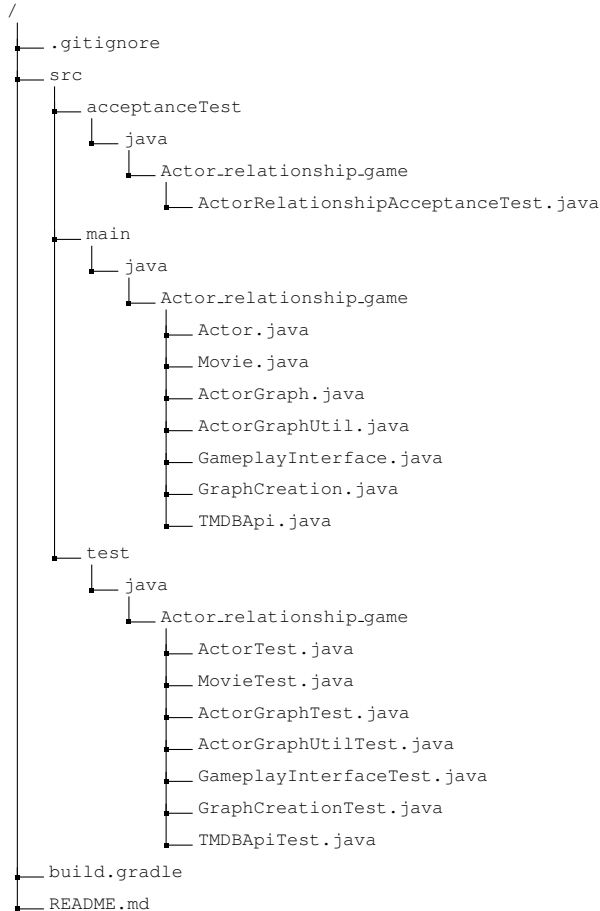
*Figure 3.* UML Class Diagram for the Example Repository



*Figure 4.* UML Sequence Diagram for the Example Repository

13

A.1.3. ARCHITECTURE DESIGN

Architecture design using file tree representation refers to a method of visualizing and organizing the structural layout of a software system's components in a hierarchical format. This approach delineates the organization of software modules, packages, libraries, and other assets in a tree-like structure, where each node represents a file or a directory containing more files or directories. Such a representation is crucial in conveying the architectural blueprint of a software project, illustrating how its various parts are interrelated.

The text-based representation of the file tree for the Actor Relationship Game repository, including test classes for each Java class, is shown as below.

```
/
├── .gitignore
├── src
│   ├── acceptanceTest
│   │   └── java
│   │       └── Actor_relationship_game
│   │           └── ActorRelationshipAcceptanceTest.java
│   ├── main
│   │   └── java
│   │       └── Actor_relationship_game
│   │           ├── Actor.java
│   │           ├── Movie.java
│   │           ├── ActorGraph.java
│   │           ├── ActorGraphUtil.java
│   │           ├── GameplayInterface.java
│   │           ├── GraphCreation.java
│   │           └── TMDBApi.java
│   └── test
│       └── java
│           └── Actor_relationship_game
│               ├── ActorTest.java
│               ├── MovieTest.java
│               ├── ActorGraphTest.java
│               ├── ActorGraphUtilTest.java
│               ├── GameplayInterfaceTest.java
│               ├── GraphCreationTest.java
│               └── TMDBApiTest.java
├── build.gradle
└── README.md
```

**A.2. Software Development**

Software development is the comprehensive process of programming, documenting, optimization, and fixing involved in creating and maintaining applications, frameworks, or other software components. It encompasses all the activities that result in software products and involves a series of steps known as the software development lifecycle (SDLC).

- **Environment Setup** is the process of preparing and configuring the necessary hardware and software tools required to build and run software applications. This setup is crucial to provide a consistent, controlled, and efficient workspace for developers to code, test, and deploy their applications. The environment can be set up on an individual's local machine, on a remote server, or in a containerized environment.

- **Implementation** is when developers write code according to the software design documents, using programming languages and tools suitable for the repository.

## A.3. Quality Assurance

Quality Assurance (QA) is the systematic process of ensuring that the software being developed meets the specified quality standards and requirements before it is released.

Software testing is an integral part of QA; involves the execution of a software component or system component to evaluate one or more properties of interest. Software testing typically includes:

- **Unit Testing** is the process of testing individual units or components of a software application to ensure their behaviors. A unit is the smallest testable part of any software and usually has one or a few inputs and usually a single output. In procedural programming, a unit could be an entire module, but it is more commonly an individual function or procedure.

- **Acceptance Testing** is a level of software testing where a system is tested for acceptability. It provides the final assurance that the software meets the PRD and is ready for use by end-users.

Code Listing 1 is part of the unit test suite for the Actor Relationship Game Repository, specifically designed to validate the functionality of the `Actor` class. Using the JUnit framework, it defines two test cases: `testActorIdAndName` and `testMovieIds`. The first test, `testActorIdAndName`, instantiates an `Actor` object with a specific ID and name ("101" and "John Doe", respectively) and asserts that the `getId()` and `getName()` methods correctly return these values, ensuring the actor's identity is accurately stored and retrievable. The second test, `testMovieIds`, creates another `Actor` object and adds two movie IDs ("201" and "202") to the actor's list of movie IDs. It then verifies that these movie IDs are indeed associated with the actor by checking if the actor's `getMovieIds()` set contains the added IDs. Together, these tests check the integrity of the `Actor` class's basic functionalities: maintaining an actor's identity and managing their associated movie IDs.

*Listing 1.* Example Unit Test

```
package Actor_relationship_game;

import org.junit.jupiter.api.Test;
import static org.junit.jupiter.api.Assertions.*;

class ActorTest {
    @Test
    void testActorIdAndName() {
        Actor actor = new Actor("101", "John Doe");
        assertEquals("101", actor.getId());
        assertEquals("John Doe", actor.getName());
    }

    @Test
    void testMovieIds() {
        Actor actor = new Actor("102", "Jane Smith");
        actor.getMovieIds().add("201");
        actor.getMovieIds().add("202");
        assertTrue(actor.getMovieIds().contains("201"));
        assertTrue(actor.getMovieIds().contains("202"));
    }
}
```

In Code Listing 2, the example acceptance test is designed to verify the functionality of generating and comparing actor lists from graph data. It employs the `runGradleTask` method to execute specific Gradle tasks for creating graph data and generating actor lists into specified file paths, using parameters for file names to differentiate between reference and test data. The test first runs the `runGraphCreation` task with paths for both reference and test graphs, followed by the `runActorGraphUtil` task to generate actor lists from these graphs into specified file paths. Once the actor lists are generated, the test reads lines from both the reference and test actor list files and then iterates through each line in the reference actor list, followed by an assertion that each actor from the reference list is also present in the test list. This process effectively checks the integrity and consistency of the actor list generation feature by ensuring that the test actor list replicates the reference list accurately, thereby validating the application's capability to process and output graph-related data correctly.

*Listing 2.* Example Acceptance Test

```
@Test
public void testActorList() throws IOException,InterruptedException{
    runGradleTask("runGraphCreation -PfileName="+referenceGraphPath);
    runGradleTask("runGraphCreation -PfileName="+testGraphPath);
    runGradleTask("runActorGraphUtil -PgraphPath="+referenceGraphPath+" -PfilePath="+referenceActorPath);
    runGradleTask("runActorGraphUtil -PgraphPath="+testGraphPath+" -PfilePath="+testActorPath);


    List<String> referenceLines = Files.readAllLines(Paths.get(referenceActorPath));
    List<String> testLines = Files.readAllLines(Paths.get(testActorPath));

    for (String referenceLine:referenceLines){
        assertTrue(containsLine(testLines,referenceLine));
    }
}
```

# B. Metrics of Software Design Evaluation

**Evaluation Guidance for UML Class**

**General Principles**

- Cohesion and Decoupling: The design should aim for high cohesion within individual classes and low coupling between different classes. High cohesion ensures that each class is dedicated to a singular task or concept, enhancing clarity and functionality. Low coupling reduces dependencies among classes, facilitating easier maintenance and scalability.
- Complexity: Utilize metrics such as the total number of classes, the average number of methods per class, and the depth of the inheritance tree to evaluate complexity. It's important to discern between conceptual classes and attributes; not every noun should become a class. The complexity level should be appropriately balanced, aligning with the specific requirements detailed in the repository's Product Requirement Document (PRD).
- Practicability: A practical design should be readable and understandable, offering a clear and comprehensive representation of the software's structures, functionalities, and behaviors. This enhances ease in programming, testing, and maintenance. Modularity should be evident, with each component serving a distinct function, streamlining the development process. Interfaces need to be designed for simplicity, facilitating smooth interactions within the software and with external environments. The design must also support robust testing strategies, enabling thorough validation through unit and acceptance tests, ensuring the design's viability in real-world applications.

**Faithfulness**

- Ensure that the design aligns with the given PRD strictly, achieving all the functionalities based on the requirements without making any hallucinations and additions. Ensure that the conceptual classes and their relationships accurately represent the essentials outlined in the PRD. This includes a detailed focus on the associations between classes, their cardinalities, and the types of relationships such as inheritance, aggregation, and composition. Clarity in class names and the optional inclusion of attributes are key for aligning with the repository's vision.

---

**Evaluation Guidance for UML Sequence**

**General Principles**

- Uniformity and Integration: The design should demonstrate a consistent style and integrated approach, ensuring all components work seamlessly together.
- Cohesion and Decoupling: Evaluate the sequence diagram for its cohesion within sequences and coupling between different parts of the system. The goal is to ensure each sequence is focused, with minimal dependencies between different system components. Strive for high cohesion within sequences and low coupling between them.
- Interaction complexity: This metric assesses the interaction complexity of the sequence diagram, focusing on the number of messages, depth of nested calls, and the number of participating objects. It also examines how the sequence of messages and the roles of key objects are portrayed in these interactions. The ideal level of complexity should be in line with the specific requirements detailed in the repository's PRD
- Practicability: This comprehensive metric includes aspects of readability, understandability, class and method representation, and the overall clarity in depicting system interactions and functionalities. Evaluate the diagram's ease of interpretation for development, testing, and maintenance, its ability to represent the functionality and purpose of each class, document object creation instances, and demonstrate the modularity and interface simplicity that support efficient and reliable system operation.

---

**Faithfulness**

- Evaluate how accurately and comprehensively the sequence diagram reflects the system's intended behavior and requirements specified in the PRD. This includes how well it captures system events, both with and without parameters, and the accuracy with which it reflects the impact of these events on the system's behavior. Also Evaluate how accurately and comprehensively the sequence diagram reflects the structural design outlined in the given UML class diagrams, ensuring a coherent and consistent development process.

---

**Evaluation Guidance for Architecture Design**

**General Principles**

- Uniformity and Integration: The design should demonstrate a consistent style and integrated approach, ensuring all components work seamlessly together, ensuring high cohesion and decoupling.
- Distinction Between Design and Coding: Recognize that the design process is distinct from coding; good design lays the groundwork for effective coding but is not synonymous with it.
- Practicability: Evaluate the architecture's practicability by assessing its organization, readability and modularity, and efficiency. The design should feature a logical and clear structure, evidenced by a well-organized file tree and distinct class locations in proper directories.
- Conformance: Evaluate the architecture for its conformance to community and industry standards. The file tree structure, coding practices including naming conventions, documentation and other structural elements should adhere to the widely accepted conventions by the open-source community and best practices of the programming language used, such as C/C++, Python, Java and JavaScript.

---

**Faithfulness**

- The architecture must be in strict accordance with the given PRD and UML class diagrams. It should accurately reflect the requirements specified in the PRD and the structural design outlined in the UML diagrams, ensuring a coherent and consistent development process.

## C. Discussions

### C.1. Model Capacity

**Challenges in Creating C/C++ Makefile and Java Gradle** The detailed results broken down on different languages are illustrated in Figure C. Makefiles and Gradle files are pivotal in C/C++ and Java repositories, orchestrating file organization, compilation, linking, and managing dependencies. However, models struggle in accurately generating these files. Smaller
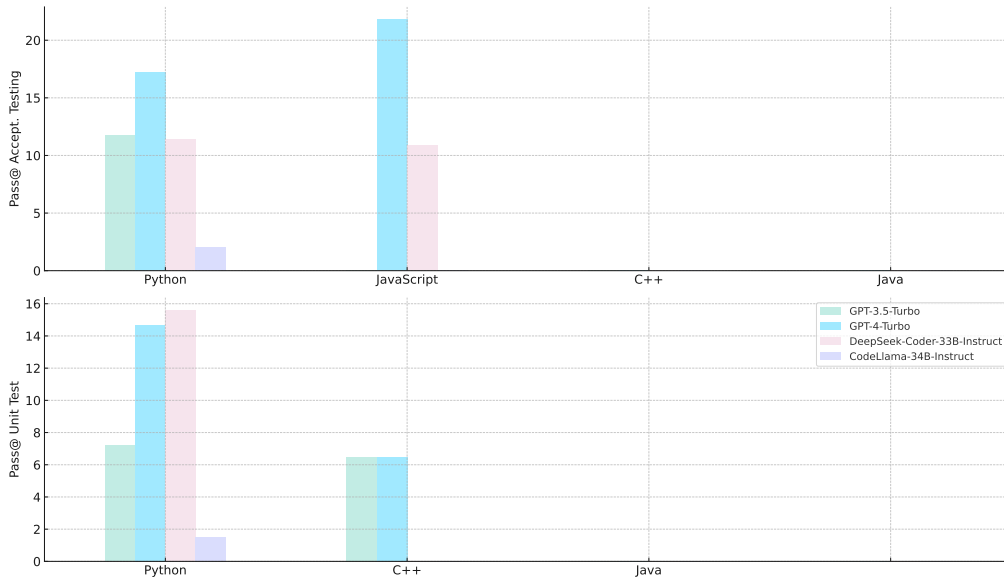
*Figure 5.* Performance Break Down on Different Languages. The results are averaged across all repositories and weighted by the number of code lines

models often misconstrue the syntax and structure, leading to irrelevant content like misplaced Linux commands or repeating given prompt instructions. Although larger models, such as GPT-4-Turbo, show improvements, they still fall short, exemplified by frequent "missing separator" errors in Makefiles. Similarly, Gradle files often lack crucial dependency declarations, a critical oversight leading to Java compilation failures. Another frequent Gradle problem is omitting essential tasks. Although given with necessary "tests" and "acceptanceTests" prompts, the model often misses these in the Gradle file, but can correct them upon review.

**Function Redefinition in Multi-file Repositories**    Models face significant challenges in multi-file repositories contexts, particularly with function redefinitions. Specifically, if a global function is required for the entire repository, it can be defined in any file, and other files just need to correctly reference it. But they tend to redundantly implement the same function across multiple files, suitable for single-file repositories but erroneous in multi-file scenarios. In C/C++, models incorrectly handle header (.h) and implementation (.cpp) files, leading to redundant declarations and conflicting implementations. These issues highlight a gap in the models' understanding of file-specific roles in programming languages.

**File Reference and Linkage Errors**    Correct file referencing is essential in multi-file programming repositories. Models, especially those with larger parameters, generally perform well in establishing basic reference logic, such as using "import" in Python and "#include" in C++. However, without review, reference errors are common, likely due to the models' sequential code generation approach, which limits their ability to correct earlier mistakes. More complex reference issues also arise. Models often struggle to differentiate between global functions and class methods, leading to reference errors when attempting to access class methods directly. These errors are challenging to rectify through review. This indicates a need for models to better grasp the intricacies of programming language structures and conventions.

## C.2. Instruction Following

**Naming Errors**    Proper naming in code is crucial for readability and maintainability. Larger models, like GPT-4-Turbo, while capable of generating syntactically correct code, exhibit deficiencies in this area. It inaccurately modifies function and variable names, disrupting the code's functionality. For instance, in "Graph BFS DFS", we instruct the use of the "top()" method for stack access, yet GPT-4-Turbo incorrectly labels it as "getTop()". Furthermore, the model's lack of attention to plurals and capitalization aggravates these errors. Despite reviews, this issue remains inadequately addressed.

**Function Parameters and Overloading Errors**    We observe that even some large models, such as GPT-4-Turbo, neglect the correct number of function parameters, missing critical ones. Function overloading is another similar nuanced aspect that many models mishandle. They often overlook the necessity of multiple constructors or methods with varying parameters. For example, in a task like "area_calculation", models fail to create both parameterized and parameterless constructors, focusing solely on the former. This oversight is not significantly rectified in the review stages, as models mistakenly attribute the errors to the test program. They stubbornly resist correcting their generated code, even when provided with clear instructions.

**Type Errors**    Models demonstrate a lack of sensitivity to type conversions, especially in strongly-typed languages like C/C++. Errors in matching const types and misusing pointer types are common, and these missteps are not readily resolved in the review process. However, in weakly-typed languages like Python, such issues are less critical but still present a concern for code accuracy.

**Variable Scope and Lifecycle Mismanagement**    Models frequently misuse variables beyond their intended scope or lifecycle. For example, they might attempt to use a loop control variable outside its loop. Another issue is the misunderstanding of private and public members in classes, where models inappropriately access private elements from outside the class. This indicates a gap in the models' understanding of encapsulation and scope management in object-oriented programming.

### C.3. Hallucination

**Fabrication of Variables**    A significant challenge is the models' propensity to fabricate non-existent local variables, a problem typically rectified during review. This issue suggests a fundamental limitation in the models' sequential generation process. Unable to retroactively integrate essential variable definitions, the models end up introducing imaginary variables in the code, resulting in apparent inaccuracies.

**Misinterpretation of Data Files**    Models also exhibit a tendency to incorrectly interpret data files as Python libraries. They attempt to import methods from these non-existent libraries, leading to further reference errors. This behavior underscores the intricacy involved in handling file references accurately within code generation tasks.

### C.4. Limitations in Testing

In Acceptance Testing task, we employ an execution-based evaluation method for the generated tests, foregoing manual quality assessments. This approach assumes a test is valid if it can accurately assess standard implementation code. However, we observed that smaller models, such as DeepSeek-Coder-6.7B-Instruct, tend to game this method. They set arbitrary criteria and invariably provide positive feedback, thereby circumventing a genuine evaluation.

Larger models like GPT-4-Turbo fell short of our expectations. They persistently recall and utilize methods in the original repository code, instead of our specially designed versions, leading to frequent import errors. This issue is exemplified in our tests with the "GeoText" and "Stocktrends" repositories. We modified the original repositories by removing "__init__.py" files, expecting models could correctly handle import relationships without them, based on our provided file structures. However, the models continued to follow the import logic of the original repositories, leading to hallucination and inaccurate test generation. This indicates a training data bias, where these models are predisposed to original repository code and show a reluctance to adjust to new circumstance.

## D. Repositories statistics in DevBench

| Text | Language | Domain | #code files | #code lines | #code tokens | #acceptance tests | #unit tests | Unit test coverage |
|---|---|---|---|---|---|---|---|---|
| TextCNN | Python | DL, NLP | 5 | 403 | 1566 | 1 | 10 | 99 |
| ArXiv digest | Python | SE, API | 1 | 198 | 901 | 4 | 38 | 94 |
| chakin | Python | NLP | 1 | 62 | 225 | 1 | 1 | 86 |
| readtime | Python | ALGO | 3 | 284 | 920 | 4 | 8 | 95 |
| hone | Python | SE | 4 | 274 | 844 | 5 | 7 | 90 |
| Stocktrends | Python | ALGO | 1 | 384 | 1350 | 2 | 7 | 85 |
| GeoText | Python | NLP | 2 | 470 | 1701 | 5 | 4 | 98 |
| lice | Python | SE | 2 | 376 | 1329 | 6 | 25 | 88 |
| PSO | Python | ALGO | 2 | 168 | 578 | 1 | 5 | 93 |
| hybrid images | Python | ALGO, CV | 1 | 144 | 746 | 1 | 19 | 90 |
| Actor Relationship Game | Java | ALGO, API | 8 | 493 | 1453 | 4 | 16 | 64.32 |
| Leftist Trees and Fibonacci Heaps Comparison | Java | ALGO | 3 | 632 | 2009 | 2 | 2 | 45.32 |
| Redis | Java | SE, DB | 9 | 779 | 2546 | 1 | 17 | 78.6 |
| idcenter | Java | SE | 4 | 333 | 1140 | 3 | 4 | 54.2 |
| image similarity | Java | SE, CV | 3 | 382 | 1397 | 2 | 2 | 71.34 |
| xlsx2csv | C/C++ | SE | 10 | 476 | 1440 | 5 | 8 | 95.17 |
| people management | C/C++ | SE, DB | 6 | 540 | 2043 | 7 | 9 | 95.14 |
| Area Calculation | C/C++ | SE | 7 | 162 | 307 | 3 | 3 | 90.48 |
| Graph BFS DFS | C/C++ | ALGO | 5 | 667 | 2828 | 5 | 22 | / |
| Logistic Management System | C/C++ | SE | 7 | 630 | 2007 | 7 | 17 | 99.11 |
| listen-now-frontend | JS | Web | 6 | 232 | 492 | 1 | 0 | / |
| register | JS | Web | 6 | 223 | 741 | 3 | 0 | / |

*Table 8.* Repositories statistics in DevBench