# Top Leaderboard Ranking = Top Coding Proficiency, Always? WEVOEVAL: Evolving Coding Benchmarks via LLM

Chunqiu Steven Xia\* Yinlin Deng\*

**Lingming Zhang** 

University of Illinois Urbana-Champaign

{chunqiu2, yinlind2, lingming}@illinois.edu

#### Abstract

LLMs have become the go-to choice for code generation tasks, with an exponential increase in the training, development, and usage of LLMs specifically for code generation. To evaluate the ability of LLMs on code, both academic and industry practitioners rely on popular handcrafted benchmarks. However, prior benchmarks contain only a very limited set of problems, both in quantity and variety. Further, due to popularity and age, many benchmarks are prone to data leakage where example solutions can be readily found on the web and thus potentially in training data. Such limitations inevitably lead us to inquire: Is the leaderboard performance on existing benchmarks reliable and comprehensive enough to measure the program synthesis ability of LLMs? To address this, we introduce 💥 EVOEVAL- a program synthesis benchmark suite created by *evolving* existing benchmarks into different targeted domains for a comprehensive evaluation of LLM coding abilities. Our study on **51** LLMs shows that compared to the high performance obtained on standard benchmarks like HUMANEVAL, there is a significant drop in performance (on average 39.4%) when using EVOEVAL. Additionally, the decrease in performance can range from 19.6% to 47.7%, leading to drastic ranking changes amongst LLMs and showing potential overfitting of existing benchmarks. Furthermore, we showcase various insights, including the brittleness of instruction-following models when encountering rewording or subtle changes as well as the importance of learning problem composition and decomposition. EVOEVAL not only provides comprehensive benchmarks, but can be used to further evolve arbitrary problems to keep up with advances and the ever-changing landscape of LLMs for code. We have open-sourced our benchmarks, tools, and complete LLM generations at https://github.com/evo-eval/evoeval

### 1 Introduction

Program synthesis [15] is widely regarded as the *holy-grail* in the field of computer science. Recently, large language models (LLMs) have become the default choice for program synthesis due to its code reasoning capabilities acquired through training on large amounts of open-source code repositories. Popular LLMs like GPT-4 [36], Claude-3 [3], and Gemini [43] have shown tremendous success in aiding developers on a wide-range of coding tasks such as code completion [10], repair [51], and test generation [12]. Furthermore, researchers and industry practitioners have designed code LLMs (e.g., DeepSeeker Coder [16], CodeLlama [40], and StarCoder [26]) using a variety of training methods designed specifically for the code domain to improve LLM code understanding.

In order to evaluate the coding abilities of LLMs, benchmarks like HUMANEVAL [10] and MBPP [4] have been handcrafted to evaluate the program synthesis task of turning natural language descriptions (e.g., docstrings) into code snippets. These code benchmarks measure

<sup>\*</sup>Contributed equally with author ordering decided by Nigiri.



Figure 1: Example problems generated in EVOEVAL through the use of targeted transformation prompts starting from a HUMANEVAL problem.

functional correctness by evaluating LLM-generated solutions against a set of limited predefined tests. Recent work [28] has further improved these benchmarks with augmented tests to rigorously evaluate the functional correctness of LLM generated code. However, apart from test inadequacy, existing popular code synthesis benchmarks have the following limitations:

- Limited amount and variety of problems. Code benchmarks are mainly constructed by human annotators manually. Due to the high manual effort required, they only contain a limited amount of problems. For example, HUMANEVAL [10] only contains 164 hand-crafted problems. Such a low amount of problems is not sufficient to fully measure the complete spectrum of program synthesis capability of state-of-the-art LLMs. Additionally, these code benchmarks include mostly self-contained coding problems that lack variety in both problem types and domains, where the final evaluation output only shows the percentage of problems solved. While they provide a baseline overview of the coding abilities, LLM builders and users cannot gain deeper insights to exactly what problem types or coding scenarios the particular LLM may excel or struggle in.
- Prone to data leakage and training dataset composition. Popular benchmarks like HU-MANEVAL and MBPP were released almost 4 years ago, with example solutions available in third-party open-source repositories. While recent LLMs have been taking turns climbing the leaderboard by achieving higher pass@1 scores (often with less than 1 percent difference between the next best model), just how much of that is attributed to having leaked solutions as part of the training data? Furthermore, the problems within these benchmarks are often simple derivatives of common coding problems/concepts. In fact, recent work [39] has shown that there are substantial overlap between benchmark solutions and open-source training corpuses. In addition, closed-source LLMs may even deliberately include benchmark groundtruths to artificially boost their leaderboard status [7]. As such, it is unclear whether high scores achieved by LLMs are truly due to their learnt coding capability or instead obtained via memorizing benchmark solutions.

As more LLMs are being constructed, trained, and used especially for code, the insufficient evaluation benchmarks raise the question of validity: *Is leaderboard performance on existing benchmarks reliable and comprehensive enough to measure the program synthesis ability of LLMs?* 

**Our work.** To address the limitation of existing benchmarks, we introduce EVOEVAL<sup>1</sup> – a set of program synthesis benchmarks created by *evolving* existing problems. The key idea behind EVOEVAL is to use LLMs instead of humans to produce new code synthesis problems based on a variety of different instructions aimed at evolving or transforming the existing benchmark problems into targeted domains for more comprehensive evaluation. Different from prior benchmark constructions that either obtain problems from open-source repositories or databases – leading to data leakage or require manual construction of each problem – resulting in high manual effort and limited diversity, EVOEVAL directly uses LLMs with targeted transformation prompts to synthesis new coding problems. Specifically, we design 5 different targeted transformation prompts: *Difficult, Creative, Subtle, Combine and Tool Use.* We then prompt GPT-4 to independently transform any existing problem in previous benchmarks into a new problem in the targeted domain.

Figure 1 shows a concrete example of EVOEVAL in action starting with an initial problem in HUMANEVAL- vowel\_counts to count the number of vowels in the string. 
We first observe the transformation to a more difficult problem by asking GPT-4 to add additional constraints or requirements. This new problem contains a separate custom vowel list that makes the overall program logic more complex. 2 We can also transform to a more creative problem of create\_alias that still uses concepts like vowels and consonants but involves a much more creative and unusual problem description. (3) We can also make subtle changes to the problem where we only count the lowercase vowels to test if the LLM is simply memorizing the benchmark. 
 We can additionally combine concepts from multiple problems together. In the example, we use another problem bf to create a new problem that returns the vowels in each planet sorted based on the orbiting order. **S** Furthermore, we can test the ability for LLMs to utilize auxiliary helper functions (common place in real-world code repositories) to solve more complex problems. Again we reuse the concepts of vowels from the initial problem, where the frequency of each vowel should be computed. However instead of directly solving the problem, the LLM can directly use the provided check\_vowel helper function to simplify the solution.

Together, each of these transformed benchmarks are designed to introduce more difficult and complex problems as well as test different aspects of the LLM code understanding and synthesis ability. In EVOEVAL, we additionally use GPT-4 to generate the groundtruth solution to each problem as well as rigorous test cases to ensure we can evaluate the functional correctness of LLM-synthesized code on EVOEVAL. Finally, we manually check each generated problem and corresponding groundtruth to ensure problem clarity and correctness. EVOEVAL serves as a way to further evolve existing benchmarks into more complex and well-suited problems for evaluation in order to keep up with the ever-growing LLM research.

**Contribution.** Our work proposes to evolve existing problems for benchmark creation:

- **Benchmark**: We present EVOEVAL– a set of program synthesis benchmarks created by evolving existing popular HUMANEVAL coding benchmark problems. EVOEVAL includes 828 problems across 5 semantic-altering and 2 semantic-preserving benchmarks. Furthermore, EVOEVAL also includes additional benchmarks to study program synthesis concepts like problem composition and decomposition. EVOEVAL is fully complete with groundtruth implementations and robust testcases to evaluate functional correctness.
- **Approach**: We propose a complete pipeline to directly synthesize new coding problems for benchmarking by evolving existing problems through the use of targeted transformation prompts. Our pipeline aims to reduce manual checking effort using a self-consistency approach to automatically refine any problem inconsistencies and generate groundtruth as well as test cases. Our approach is general and can be used on other benchmark problems, adopted for transformation into additional domains or utilize different problem generation strategies [50].
- **Study**: We conduct a comprehensive study on **51** different LLMs across all benchmarks in EVOEVAL. We found that compared to the high performance obtained on standard benchmarks like HUMANEVAL, when evaluated on EVOEVAL, popular LLMs significantly drop in performance (on average 39.4%). Additionally, this drop is not uniform across all

<sup>&</sup>lt;sup>1</sup>coincidentally similar pronunciation with EVILEVAL



Figure 2: Overview of EVOEVAL evolving problem generation pipeline.

LLMs and can range from 19.6% to 47.7%, leading to drastic ranking changes amongst top performing models. We further demonstrate that certain LLMs cannot keep up their high performance obtained in HUMANEVAL when evaluated on more challenging or problems in different domains, highlighting the possibilities of overfitting to existing benchmarks. Moreover, we observe that while instruction-following LLMs perform well in solving self-contained problems, they struggle with the tool using aspect of utilizing already provided auxiliary functions. Furthermore, they are particularly sensitive to the problem description where rephrasing or subtle changes to the problem docstring leads to degradation in output solutions compared to their base non-instruction-following counterparts. Additionally, we demonstrate that current state-of-the-art LLMs fail to effectively compose multiple general coding concepts to solve more complex variants, or address subproblems decomposed from previously solved difficult problem.

## 2 Approach

Figure 2 shows the overview of the benchmark creation pipeline for EVOEVAL. We start by taking the original problem and apply a chosen targeted transformation prompt aimed at prompting GPT-4 to produce a new code synthesis problem along the targeted domain. Using this initial transformed problem, we enter our refinement pipeline to fix any ambiguities or inconsistencies in the problem description, as well as generating the test cases and groundtruth solution for functional evaluation. Finally, to ensure correctness, we manually examine each produced problem along with the groundtruth and make corresponding changes to produce the final evolved benchmarks.

**Targeted problem transformation.** EVOEVAL uses zero-shot prompting to evolve an existing coding benchmark to produce new and diverse problems. Each transformation prompt, as shown in the examples in Figure 1, aims to transform the existing problem in a specific manner. In particular, we define two different types of transformation prompts: 1) **semanticaltering** – change the semantic meaning of the original problem and 2) **semantic-preserving** – modify the problem description while keeping the semantic meaning the same. While Figure 1 shows only semantic-altering transformation prompts to produce new problems, we can also produce semantic-preserving problems to test additional aspect of the LLM coding abilities.

**Problem refinement & groundtruth Generation.** The initial evolved problem produced by GPT-4 may include small inconsistencies such as contradicting sentences or incorrect I/O examples in the docstring. For coding benchmarks, such inconsistencies are especially damaging as it can detract from the problem specification, leading to inaccurate evaluation of LLM coding capabilities. As such, we introduce a refinement pipeline to iteratively rephrase and refine problem as needed. In addition, during this process, we also use GPT-4

Table 1: EVOEVAL and HUMANEVAL benchmark statistics. Note: the number in bracket shows the number of testcases in the augmented HUMANEVAL+ benchmarks, in EVOEVAL, they are directly reused in SUBTLE, VERBOSE and CONCISE due the similarity.

|                   | original    |           | set      | semantic-preserving |         |          |             |             |
|-------------------|-------------|-----------|----------|---------------------|---------|----------|-------------|-------------|
|                   | HUMANEVAL   | DIFFICULT | CREATIVE | SUBTLE              | Combine | TOOL_USE | VERBOSE     | CONCISE     |
| # problems        | 164         | 100       | 100      | 100                 | 100     | 100      | 164         | 164         |
| Avg. problem len. | 450.6       | 749.4     | 982.1    | 406.8               | 860.4   | 1224.6   | 450.6       | 450.6       |
| Avg. # test cases | 9.6 (764.1) | 49.8      | 43.1     | 10.3 (745.4)        | 51.8    | 51.3     | 9.6 (764.1) | 9.6 (764.1) |



(a) CREATIVE & TOOL\_USE (b) SUBTLE, DIFFICULT, COMBINE (c) VERBOSE & CONCISE

Figure 3: 2 dimensional t-SNE visualization of EVOEVAL benchmarks.

to produce the necessary groundtruth implementation of the function as well as example test cases to be used for evaluation.

We first directly use GPT-4 to obtain a possible solution for the initial problem. Additionally, we also prompt GPT-4 to extract (if available in the initial problem docstring) or produce the test inputs for the transformed problem. We then evaluate the test inputs on the solution to derive the corresponding expected test outputs. Next, using these test inputs/outputs, we instruct GPT-4 to add or fix the example test cases in the docstring, providing further demonstrations of the task.

Using this refined problem, we again generate a solution. We then leverage selfconsistency [47] to check if the new solution on the test inputs produce the same outputs as the previous solution. The intuition is that since both solutions are generated by GPT-4 and the refined problem should only include minimal changes (e.g., adding new testcase examples), the solution output should then be the same in the absence of any potential inconsistencies or ambiguity in problem description. As such, if we observe differences between the two solution outputs, we ask GPT-4 to further rephrase and fix any inconsistencies in the original problem and repeat the process. On the other hand, if both solutions agree on outputs, we terminate the problem refinement stage and return the trio comprising of the new problem description, the solution as the groundtruth and the test cases for functional evaluation.

**Manual examination & test augmentation.** For each transformed problem, we carefully examine and adjust any final faults to ensure each problem and groundtruth is correctly specified and implemented. Additionally, using the initial set of test cases from the refinement stage, we further generate additional tests following the LLM-based test augmentation technique in EVALPLUS [28]. Finally, we produce EVOEVAL, a comprehensive code synthesis benchmark suite, which through the use of evolving transformations can generate diverse coding problems to evaluate LLM coding capability across various problem domains.

### 3 EVOEVAL Dataset Overview

We use the problems in HUMANEVAL as seeds to produce EVOEVAL. Problems in EVOEVAL consist mainly of self-contained functions, except for TOOL\_USE that includes helper functions specifically designed to test the tool using capability of LLMs. Each problem uses a

docstring to illustrate the problem specification, along with test cases and groundtruth to evaluate the functional correctness. Table 1 shows the statistics of the benchmarks in EVOE-VAL. In total, EVOEVAL includes 828 problems across 7 different datasets (5 semantic-altering and 2 semantic-preserving):

- **DIFFICULT**: Introduce complexity by adding additional constraints and requirements, replace commonly used requirements to less common ones, or add additional reasoning steps to the original problem.
- CREATIVE: Generate a more creative problem compared to the original through the use of stories or uncommon narratives.
- **SUBTLE**: Make a subtle and minor change to the original problem such as inverting or replacing a requirement.
- **COMBINE**: Combine two different problems by integrating the concepts from both problems. In order to select problems that make sense to combine, we apply a simple heuristic to combine only problems of the same type together categorized based on the type of input arguments in the original problem.
- **TOOL\_USE**: Produce a new problem containing a main problem and one or more helpers functions which can be used to solve it. Each helper function is fully implemented and provides hints or useful functionality for solving the main problem. The main problem does not explicitly reference individual helper functions, and we do not require the model to use the provided helpers.
- VERBOSE: Reword the original docstring to be more verbose. These verbose docstrings can use more descriptive language to illustrate the problem, include detailed explanation of the example output, and provide additional hints.
- **CONCISE**: Reword the original docstring to be more concise by removing unnecessary details and using concise language. Furthermore, simple examples that are not required to demonstrate edge cases may be removed.

For each of the semantic-altering benchmarks, we generate 100 problems each using different seed problems from HUMANEVAL. For semantic-preserving benchmarks, we generate using all 164 problems in HUMANEVAL as it requires less validation since we can reuse the original groundtruths. As shown in Table 1, compared to HUMANEVAL, EVOEVAL contains longer coding questions with longer average problem length. Furthermore, EVOEVAL also uses more test cases to perform robust evaluation compared to base HUMANEVAL.

Figure 3 shows the embedding visualization using t-SNE [18]<sup>2</sup> by projecting high-dimension representation of the problems docstrings in both EVOEVAL and HUMANEVAL into the 2D plane. First, we see that CREATIVE and TOOL\_USE drastically change the embedding distribution compared to the original dataset. The arrow in Figure 3a shows one example of the shift in distribution from the original problem to a creative one. Next, we see that SUBTLE, DIFFICULT and COMBINE largely retain the same distribution as the original problems. This is due to the high parity across these problem descriptions where SUBTLE only applies subtle changes and DIFFICULT adds additional complex constraints while keeping the main problem descriptions largely the same. Specifically, for COMBINE, we can see from an example arrow in Figure 3b, the new combined problem shifts the embedding for both of the original problems. Finally, we observe that for VERBOSE and CONCISE, the embeddings almost perfectly match the original problem, reflecting their semantic-preserving nature. In Appendix C, we present example problems for each benchmark in EVOEVAL.

## 4 Methodology

**Setup.** Each LLM generated sample is executed against the test cases in EVOEVAL and evaluated using differential testing [31] – comparing against the groundtruth results to measure functional correctness. We report the functional correctness by using the popular pass@k metric. We focus on greedy decoding (i.e., producing a deterministic sample per each problem with temperature = 0). We denote this as pass@1.

<sup>&</sup>lt;sup>2</sup>perplexity=50 and iter=1000 using text-embedding-3-large model from OpenAI



Figure 4: Comparison of pass@1 on HUMANEVAL+ and EVOEVAL datasets of selected models. The red dotted identity line (i.e., x = y) represents equivalent performance on both HUMANEVAL and EVOEVAL. For each benchmark, we cluster the LLMs into 1) purple region – aligned performance on HUMANEVAL and EVOEVAL and 2) blue region – over performant LLMs on HUMANEVAL compared with EVOEVAL results.



Figure 5: Ranking changes across EVOEVAL benchmarks of (a) top-10 and (b) bottom-10 best performing LLMs on HUMANEVAL respectively. Y-axis shows the normalized pass@1 score defined as the LLM pass@1 normalized by the minimal and maximum pass@1 achieved by all LLMs on benchmark.

**Models.** We evaluate **51** popular state-of-the-art LLMs, including both proprietary and open-source models on EVOEVAL. We evaluate not only the popular general-purpose LLMs but also include recent code-based LLMs for comprehensive evaluation. Further, we classify the LLMs as either base or instruction-following and focus our analysis on discussing the effect of model variants have on EVOEVAL performance.

**Input format.** To produce the code solution using each LLM, we provide a specific input prompt: For base LLMs (i.e., not instruction-tuned variants), we simply use only the function header with the docstring and let the LLM autocomplete the solution. For instruction-following LLMs, we follow the model-makers' guide on the exact instruction and format to use and ask the LLM to generate a complete solution for the problem.

### 5 Evaluation

### 5.1 LLM Synthesis & Evaluation on EVOEVAL

**EVOEVAL produces more complex and challenging benchmarks for program synthesis.** Table 2 shows the pass@1 performance along with the ranking of LLMs on each of the semantic-altering EVOEVAL benchmarks with the average pass@1 and ranking on all benchmarks in the last columns. First, compared to the success rate on HUMANEVAL, when evaluated on EVOEVAL, all LLMs **consistently perform worse**. For example, the state-of-the-art GPT-4, GPT-4-Turbo and Claude-3 models solve close to 85% of all HU-MANEVAL problems but fall almost below 50% pass@1 when evaluated on the DIFFICULT problems. On average, across all benchmarks, the performance of LLMs decreased by 39.4%

|                                | Size HUMANE                          | VAL        | ₩Diffi       | CULT       | <b>C</b> REA | TIVE       | ₩SUB           | ГLЕ        | ₩Сом        | BINE       | ₩Tooi        | _USE       | ₩EvoE        | VAL             |
|--------------------------------|--------------------------------------|------------|--------------|------------|--------------|------------|----------------|------------|-------------|------------|--------------|------------|--------------|-----------------|
|                                | pass@1                               | rank       | pass@1       | rank       | pass@1       | rank       | pass@1         | rank       | pass@1      | rank       | pass@1       | rank       | pass@1       | rank            |
| © GPT-4-Turbo⊨                 | NA 83.5 (80.5)                       | 1          | 50.0         | *2         | 61.0         | 2          | 82.0           | 1          | 45.0        | 2          | 69.0         | *1         | 65.1         | 2               |
| \$ GPT-4₽                      | NA 82.3 (76.2)                       | *2         | 52.0         | 1          | 66.0         | 1          | 76.0           | 3          | 53.0        | 1          | 68.0         | 3          | 66.2         | 1               |
| © ChatGPT ⊨                    | NA 76.8 (69.5)                       | *5         | 33.0         | *13        | 42.0         | *7         | 70.0           | 4          | 33.0        | 4          | 64.0         | *6         | 53.1         | 6               |
| Claude-3                       | NA 82.3 (75.0)                       | *2         | 50.0         | *2         | 53.0         | 3          | 81.0           | 2          | 42.0        | 3          | 69.0         | *1         | 62.9         | 3               |
| Claude-3-haiku                 | NA 74.4 (66.5)                       | *8         | 40.0         | *6         | 47.0         | *5         | 65.0           | *10        | 25.0        | *6         | 61.0         | *10        | 52.1         | 7               |
| Claude-2                       | NA 66.5 (62.2)                       | *18        | 29.0         | 17         | 42.0         | *7         | 64.0           | *13        | 19.0        | 14         | 57.0         | *16        | 46.2         | 15              |
| ● Gemini <sup>®</sup>          | NA 62.2 (56.7)                       | 21         | 37.0         | *10        | 40.0         | 12         | 53.0           | *21        | 23.0        | *9         | 57.0         | *16        | 45.4         | 17              |
| OPaLM-2ഈ                       | NA 40.2 (36.6)                       | 38         | 18.0         | *32        | 22.0         | 33         | 36.0           | *42        | 3.0         | *39        | 46.0         | *29        | 27.5         | 37              |
|                                | 33b 78.0 (73.2)                      | 4          | 47.0         | 5          | 47.0         | *5         | 67.0           | *5         | 31.0        | 5          | 66.0         | 4          | 56.0         | 4               |
| ♂ DeepSeeker-Inst <sup>™</sup> | 6.7b 74.4 (69.5)<br>1.3b 63.4 (60.4) | -78<br>20  | 40.0<br>20.0 | *6<br>*30  | 37.0<br>25.0 | *13<br>*25 | 61.0<br>53.0   | *17<br>*21 | 18.0<br>9.0 | *15<br>*28 | 51.0<br>39.0 | 24<br>*41  | 46.9<br>34.9 | 14<br>24        |
|                                | 33b 50.6 (42.7)                      | 26         | 26.0         | 20         | 23.0         | *30        | 47.0           | *26        | 11.0        | *25        | 63.0         | *8         | 36.8         | 23              |
| ♂ DeepSeeker                   | 6.7b 45.1 (38.4)                     | *31        | 21.0         | *26<br>*49 | 24.0         | *27<br>*25 | 47.0           | *26        | 5.0         | *35        | 55.0         | *19        | 32.9         | 29<br>45        |
| DoopSookar-1 5-Inct            | 7b 68 9 (63 4)                       | *15        | 37.0         | *10        | 37.0         | *13        | 66.0           | *9         | 24.0        | 8          | 40.0         | *12        | 48.8         | 10              |
| DeepSeeker-1.5-Ilist           | 7b 42 1 (24 8)                       | *25        | 21.0         | *26        | 37.0         | *17        | 42.0           | *21        | 4.0         | *27        | 54.0         | *21        | 22.0         | 28              |
| Deepseeker-1.5                 | 70 42.1 (34.8)                       | *19        | 21.0         | 16         | 41.0         | *10        | 45.0           | *10        | 18.0        | *15        | 65.0         | - 21       | 47.7         | 12              |
| @Codel lama-Inst®              | 34b 51.8 (43.9)                      | 25         | 22.0         | *24        | 27.0         | 23         | 43.0           | *31        | 9.0         | *28        | 47.0         | *27        | 33.3         | 27              |
| V Couclianta matp              | 13b 48.8 (42.7)<br>7b 43 3 (39.0)    | 29<br>33   | 21.0<br>14.0 | *26        | 25.0<br>18.0 | *25<br>*37 | $46.0 \\ 40.0$ | 29<br>*37  | 8.0<br>8.0  | *31<br>*31 | 54.0<br>44.0 | *21<br>*33 | 33.8<br>27.9 | 26<br>36        |
|                                | 70b 60 4 (52 4)                      | 23         | 25.0         | 21         | 29.0         | *20        | 49.0           | *23        | 14.0        | *21        | 63.0         | *8         | 40.1         | 22              |
| @ CodeLlama                    | 34b 52.4 (43.3)                      | 24         | 15.0         | 37         | 24.0         | *27        | 47.0           | *26        | 11.0        | *25        | 44.0         | *33        | 32.2         | 30              |
|                                | 13b 42.7 (36.6)<br>7b 39.6 (36.6)    | 34<br>39   | 18.0<br>10.0 | *32<br>*42 | 24.0<br>15.0 | *27<br>41  | 38.0<br>42.0   | *39<br>34  | 6.0<br>3.0  | 34<br>*39  | 48.0<br>44.0 | *25<br>*33 | 29.4<br>25.6 | 33<br>38        |
| WizardCoder                    | 34b 61.6 (54.3)                      | 22         | 24.0         | 22         | 32.0         | 19         | 55.0           | 20         | 17.0        | *18        | 55.0         | *19        | 40.8         | 20              |
| WizardCoder-1.1                | 33b 73.8 (69.5)                      | 10         | 48.0         | 4          | 48.0         | 4          | 66.0           | *8         | 20.0        | 13         | 64.0         | *6         | 53.3         | 5               |
| XwinCoder                      | 34b 68.9 (62.2)                      | *15        | 33.0         | *13        | 42.0         | *7         | 67.0           | *5         | 15.0        | 20         | 60.0         | *12        | 47.7         | 13              |
| Phind-CodeLlama-2              | 34b 70.7 (66.5)                      | 13         | 22.0         | *24        | 35.0         | 16         | 63.0           | 15         | 25.0        | *6         | 58.0         | 15         | 45.6         | 16              |
| Code Millenials                | 34b 73.2 (69.5)                      | 11         | 35.0         | 12         | 41.0         | *10        | 65.0           | *10        | 17.0        | *18        | 56.0         | 18         | 47.9         | 11              |
| Speechless-CL®                 | 34b 75.0 (69.5)                      | 7          | 38.0         | 9          | 37.0         | *13        | 64.0           | *13        | 23.0        | *9         | 59.0         | 14         | 49.3         | 9               |
| Magicoder-s-DS                 | 6.7b 76.8 (70.7)                     | *5         | 40.0         | *6         | 34.0         | *17        | 67.0           | *5         | 21.0        | *11        | 61.0         | *10        | 50.0         | 8               |
| Magicoder-s-CL®                | 7b 70.1 (65.9)                       | 14         | 27.0         | 19         | 26.0         | 24         | 58.0           | 19         | 11.0        | *25        | 52.0         | 23         | 40.7         | 21              |
|                                | 15b 45.1 (36.0)                      | *31        | 16.0         | *35        | 19.0         | *35        | 41.0           | *35        | 5.0         | *35        | 48.0         | *25        | 29.0         | 35              |
| StarCoder2                     | 7b 34.8 (31.1)<br>3b 31.1 (26.2)     | *40<br>44  | 12.0<br>8.0  | *39<br>*45 | 17.0<br>14.0 | 39<br>*42  | 38.0<br>31.0   | *39<br>*44 | 2.0<br>2.0  | *45<br>*45 | 46.0<br>35.0 | *29<br>45  | 25.0<br>20.2 | $\frac{40}{46}$ |
| StarCoder                      | 15b 34.8 (30.5)                      | *40        | 12.0         | *39        | 11.0         | 47         | 37.0           | 41         | 2.0         | *45        | 44.0         | *33        | 23.5         | 41              |
| Mixtral-Inst 🗩                 | 8x7b 42.1 (38.4)                     | *35        | 21.0         | *26        | 18.0         | *37        | 41.0           | *35        | 9.0         | *28        | 45.0         | *31        | 29.3         | 34              |
| Mistral-Inst-v02               | 7b 28.0 (23.2)                       | *48        | 8.0          | *45        | 16.0         | 40         | 25.0           | 50         | 3.0         | *39        | 8.0          | 51         | 14.7         | 50              |
| Mistral-Inst🦻                  | 7b 28.7 (24.4)                       | 47         | 6.0          | *48        | 8.0          | 50         | 31.0           | *44        | 3.0         | *39        | 29.0         | 48         | 17.6         | 49              |
| H Mistral                      | 7b 28.0 (23.8)                       | *48        | 8.0          | *45        | 14.0         | *42        | 30.0           | 47         | 3.0         | *39        | 38.0         | 43         | 20.2         | 47              |
| OpenChat®                      | 7b 71.3 (66.5)                       | 12         | 33.0         | *13        | 29.0         | *20        | 62.0           | 16         | 14.0        | *21        | 43.0         | 38         | 42.1         | 18              |
| stable-code                    | 3b 29.3 (25.6)                       | 46         | 10.0         | *42        | 10.0         | *48        | 31.0           | *44        | 3.0         | *39        | 41.0         | 39         | 20.7         | 43              |
| 🛇 Gemma-Inst 🖻                 | 7b 28.0 (23.2)                       | *48        | 6.0          | *48        | 10.0         | *48        | 29.0           | 48         | 2.0         | *45        | 31.0         | 47         | 17.7         | 48              |
| • Gemma                        | 7b 31.7 (25.0)<br>2b 22.0 (17.1)     | 43<br>51   | 12.0<br>2.0  | *39<br>51  | 13.0<br>6.0  | *44<br>51  | 40.0<br>24.0   | *37<br>51  | 2.0<br>2.0  | *45<br>*45 | 39.0<br>21.0 | *41<br>50  | 23.0<br>12.8 | 42<br>51        |
| Phi-2                          | 2.7b 50.0 (45.1)                     | *27        | 18.0         | *32        | 23.0         | *30        | 49.0           | *23        | 14.0        | *21        | 37.0         | 44         | 31.8         | 32              |
|                                | 72b 67.1 (61.6)                      | 17         | 28.0         | 18         | 28.0         | 22         | 61.0           | *17        | 21.0        | *11        | 47.0         | *27        | 42.0         | 19              |
| Qwen-1.5                       | 14b 50.0 (45.7)<br>7b 42 1 (37.8)    | *27<br>*35 | 20.0         | *30<br>*35 | 23.0         | *30<br>*44 | 48.0           | 25<br>*31  | 18.0        | *15        | 44.0         | *33<br>16  | 33.8         | 25<br>30        |
|                                | 14b 46 3 (13 0)                      | 30         | 23.0         | 22         | 20.0         | 3/         | 45.0           | 30         | 13.0        | 24         | 45.0         | +21        | 20.0         | 31              |
| Qwen®                          | 7b 34.1 (29.9)                       | 42         | 9.0          | 44         | 12.0         | 46         | 36.0           | *42        | 4.0         | *37        | 28.0         | 49         | 20.5         | 44              |

Table 2: pass@1 and ranking results (\* indicates tie) on the semantically-altering EVOEVAL and HUMANEVAL benchmarks. Note: Denotes instruction-following LLMs.



Figure 6: Radar graph of selected models with similar HUMANEVAL scores.

(DIFFICULT: 58.7%, CREATIVE: 50.2%, SUBTLE: 5.0%, COMBINE: 78.1%, and TOOL\_USE: 4.9%). Additionally, this drop is not uniform across all LLMs and can range from 19.6% to 47.7%.

LLMs struggles on EVOEVAL benchmarks compare to high performance achieved on HUMANEVAL. One surprising finding is that, on SUBTLE, where only small changes are made to original problem with the roughly the same level of difficulty, the average performance of LLMs drops by 24.0% across the same 100 problems. It is important to note that, as the pass@l score is generally higher on the first 100 problems than the complete 164 HUMANEVAL problems, this back-to-back performance drop is much higher than the performance drop from HUMANEVAL to SUBTLE mentioned above (which is 5.0%). Furthermore, we can also identify LLMs which struggle heavily on specific types of problems compared to their relative performance on HUMANEVAL. Figure 4 shows scatter plot of HUMANEVAL+ and EVOEVAL scores of selected LLMs. As we saw before, the significant portions of the models tends to be worse on EVOEVAL than HUMANEVAL (i.e., purple shaded region). However, there exists LLMs that have a *much* higher HUMANEVAL score compared to their performance on EVOEVAL (i.e., blue shaded region). This highlights potential data leakage of popular benchmarks where LLM performances are artificially inflated but do not translate to more difficult or other program synthesis problems.

Significant ranking changes of LLMs across different EVOEVAL benchmarks. In Figure 5, compared to the existing parity – where top models all perform similarly on HUMANEVAL, we observe drastic differences in ranking changes on EVOEVAL. We observe that while the relative difference between the top 5 models on HUMANEVAL is less than 10%, the difference on EVOEVAL on average is over 20%. Due to such saturation in top model performance, existing benchmarks may not reliably rank the program synthesis ability of each model. Taking a closer look at specific models, while Claude-3 and GPT-4 are tied for the 2nd best HUMANEVAL score, they both excel at different types of problems: GPT-4 performs best on difficult and creative problems while Claude-3 can better reason about helper functions in TOOL\_USE and are less affected by subtle changes from original HUMANEVAL. Furthermore, while GPT-4-Turbo achieves the top HUMANEVAL and HUMANEVAL+ score, it falls off compare to the base GPT-4 variant where it is worse on DIFFICULT, CREATIVE and COMBINE problems. Such evaluation cannot be gained through naively reporting existing coding benchmark performance. Overall, by evolving the original benchmark into more difficult and diverse problems of different types, EVOEVAL can provide a more holistic evaluation and ranking of the coding ability of LLMs.

**EVOEVAL can be used to comprehensively compare multiple models.** Figure 6 shows two radar graphs of two sets of LLMs. In Figure 6a, while both WizardCoder-1.1 and Phind-CodeLlama-2 are top performing LLMs and have similar HUMANEVAL scores, they perform drastically differently across the benchmarks in EVOEVAL. WizardCoder-1.1 is better on DIFFICULT and CREATIVE and Phind-CodeLlama-2 are better on COMBINE problems. This can be partially explained through the training dataset used in each LLM where WizardCoder-1.1 uses an evolving dataset to generate more complex and difficult problems



Figure 7: HUMANEVAL pass@1 separated into instruction-following and non-instruction-following LLMs with relative decrease or increase in pass@1 on VERBOSE and COMBINE.

whereas Phind-CodeLlama-2 is fine-tuned on high quality programming problems that seems to boost the ability to solve programs which combines multiple smaller programming concepts. Similar phenomenon can also be observed in Figure 6b. Different from just reporting a singular pass@k score, EVOEVAL also allows detailed analysis across the different dimension of coding capability to identify particular domains or type of synthesis questions the LLM struggles or excels in.

Instruction-following LLMs are sensitive to subtle or rephrasing of problem docstring. Unlike the semantic-altering benchmarks in EVOEVAL, the semantic-preserving problems do not always lead to a decrease in performance. Figure 7 shows the HUMANEVAL score (bar) and the relative performance drop or improvement (arrows) on VERBOSE and CONCISE separated into instruction-following and base LLMs. We observe that almost all instructionfollowing LLMs drops in performance (on average 3.4% and 4.0% decrease on VERBOSE and CONCISE respectively) when evaluated on the two semantic-preserving dataset compared to the original HUMANEVAL. This is drastically different from the non-instruction-following variants where we even observe performance improvements (on average 0.5% and 2.1% increase on VERBOSE and CONCISE respectively). VERBOSE and CONCISE do not change the semantic meaning of the original problem except reword it in either a more verbose or concise manner. Prior work [11] has shown that by smartly rephrasing the original problem description, one can further boost LLM performance and we observe the similar phenomenon here mostly only for non-instruction-following models. This further points to possibility of overfitting to the exact descriptions utilized in HUMANEVAL especially for instruction-tuned LLMs.

Additionally, even on the semantic-altering benchmark of SUBTLE, where only subtle changes to the original problem are applied, on average, instruction-following LLMs drops by 7.6% whereas base models only decreases by less than 1% relative to their HUMANEVAL performance. These findings across LLM types show that while instruction-tuning is expected to align better with detailed task instructions, it fails to distinguish between these subtle changes in docstring, indicating potential memorization or contamination of prior evaluation benchmarks.

#### 5.2 Problem Composition

**Composition problems.** The ability to compose different known concepts to solve new problems is known as *compositional generalization* [24]. This skill is essential for code synthesis, especially for complex problems in real-world programs. However, measuring compositional generalization in LLM presents a fundamental challenge since it requires controlling the relationship between training and test distributions [41]. While it is not easy to control the pre-training data of LLMs, we have more control in the testing phase. Hence, we focus on program concepts that have been demonstrated to fall within the capabilities of an LLM, and explore whether this proficiency extends to the combination of program concepts. As such, we start by taking a deeper look at the COMBINE problems evolved from combining previous HUMANEVAL problems.

Table 3: Detailed results of top performing LLMs on COMBINE and COMBINE-NAIVE. "HUMANEVAL" is categorized into "*pass both*", "*pass one*" and "*pass none*", depending on the success on the two parent problems used to create COMBINE and COMBINE-NAIVE. "COMBINE (Solved)" and "COMBINE-NAIVE (Solved)" then show the distribution of successfully solved problems in the composition dataset that came from the previous categories. "*Composition Percentage*" is defined as the percentage of "*pass both*" problems the LLM can *still* solve when combining both of these problems.

|                                    | Size | Н         | IumanEv  | AL        | Cor       | Composition |           |            |
|------------------------------------|------|-----------|----------|-----------|-----------|-------------|-----------|------------|
|                                    |      | pass both | pass one | pass none | pass both | pass one    | pass none | Percentage |
| © GPT-4⊨                           | NA   | 93        | 7        | 0         | 50        | 3           | 0         | 53.8%      |
| ₲ GPT-4-Turbo∞                     | NA   | 79        | 19       | 2         | 38        | 6           | 1         | 48.1%      |
| Claude-3 <sup>®</sup>              | NA   | 81        | 19       | 0         | 35        | 7           | 0         | 43.2%      |
| ₲ ChatGPT ়                        | NA   | 65        | 34       | 1         | 24        | 9           | 0         | 36.9%      |
| 🔮 DeepSeeker-Inst폗                 | 33b  | 71        | 27       | 2         | 29        | 2           | 0         | 40.8%      |
| 🗖 Claude-3-haiku 👳                 | NA   | 63        | 34       | 3         | 19        | 6           | 0         | 30.2%      |
| 𝔹 DeepSeeker-1.5-Inst <sup>∞</sup> | 7b   | 62        | 37       | 1         | 18        | 6           | 0         | 29.0%      |
| 📀 Gemini 🦻                         | NA   | 46        | 45       | 9         | 19        | 2           | 2         | 41.3%      |
|                                    |      | Н         | IumanEv  | AL        | Сомві     | NE-NAIVE    | (Solved)  |            |
| © GPT-4₽                           | NA   | 1018      | 55       | 1         | 766       | 7           | 0         | 75.2%      |
| ₲ GPT-4-Turbo₽                     | NA   | 863       | 195      | 16        | 407       | 61          | 3         | 47.2%      |
| Claude-3                           | NA   | 796       | 268      | 10        | 359       | 96          | 1         | 45.1%      |
| ₲ ChatGPT ₪                        | NA   | 799       | 261      | 14        | 474       | 79          | 1         | 59.3%      |
| 😻 DeepSeeker-Inst📟                 | 33b  | 740       | 304      | 30        | 462       | 95          | 5         | 62.4%      |
| 🛯 Claude-3-haiku 🖻                 | NA   | 592       | 409      | 73        | 286       | 133         | 17        | 48.3%      |
| 𝔮 DeepSeeker-1.5-Inst <sup>∞</sup> | 7b   | 634       | 372      | 68        | 393       | 130         | 17        | 62.0%      |
| 🗿 Gemini 🖻                         | NA   | 364       | 535      | 175       | 225       | 205         | 37        | 61.8%      |

First half of Table 3 shows the detailed breakdown of the COMBINE dataset results on the top 8 performing LLMs. We observe that almost all problems solved in COMBINE came from the pass both category, which is intuitive as we do not expect LLMs to solve a problem composed of subproblems that it cannot already solve. However, we see that overall, the composition percentage is quite low as only GPT-4 is able achieve greater than half. This demonstrates, for the first time, that while state-of-the-art LLMs can achieve a high pass rate on simple programming tasks in general-purpose languages like Python, they still struggle with generalizing and composing these known concepts to address more complex problems.

Naive combination problems. Since COMBINE problems are not guaranteed to not contain additional new logic or concepts, we build a simplified dataset for sequential composition. Let A and B be two separate problems with *x* as input(s) for *A*, we aim to create a new problem C with same inputs where the solution can be written as B(A(x)). To accomplish this, the new problem includes a sequential docstring by attaching the docstring of problem A followed by B. Directly concatenating them will lead to unclear descriptions, as such, for each problem in HUMANEVAL, we manually create two separate variants based on which order the problem may come in the new docstring. Figure 8 shows an example naive combination problem with the manual sequential instruction highlighted in red. Using these modified problem docstrings, we build a sequential combina-



Figure 8: COMBINE-NAIVE problem

tion dataset – COMBINE-NAIVE, containing 1074 problems by randomly combining problems filtering for input output matching (i.e., type of A(x) should equal to type of y in B(y))

The latter half of Table 3 shows the results on COMBINE-NAIVE following the same setup as COMBINE. We observe that while the composition percentage on the naive dataset improves significantly compared to the evolved COMBINE dataset, it still fails to reach near perfection, with the best LLM being able to only solve 3/4 of prior pass both problems. While existing

Table 4: Detailed results of top performing LLMs on DECOMPOSE. "HUMANEVAL" shows the pass/fail breakdown of the 50 seed HUMANEVAL problems. Each of these 50 problems, initially pass or failed, is decomposed into two subproblems. These are further categorized into "*pass both*", "*pass one*" and "*pass none*", based on whether the LLM can solve both subproblems. "*Decomp.* %" is the percentage of originally passing problems for which the LLM can solve both decomposed subproblems. Similarly, "*Recomp.* %" is the percentage of originally failing problems for which the LLM can solve both decomposed subproblems.

|                                    | Size | HUM  | anEval |                |          |           | Decomp.   | Recomp.  |           |       |              |
|------------------------------------|------|------|--------|----------------|----------|-----------|-----------|----------|-----------|-------|--------------|
|                                    |      |      |        | HUMANEVAL pass |          |           | Hum       | ianEval  | %         | %     |              |
|                                    |      | pass | fail   | both pass      | one pass | both fail | both pass | one pass | both fail |       |              |
| © GPT-4₽                           | NA   | 47   | 3      | 37             | 10       | 0         | 0         | 3        | 0         | 78.7% | 0.0%         |
| © GPT-4-Turbo∞                     | NA   | 39   | 11     | 29             | 9        | 1         | 4         | 6        | 1         | 74.4% | 36.4%        |
| Claude-3                           | NA   | 39   | 11     | 26             | 11       | 2         | 6         | 5        | 0         | 66.7% | 54.5%        |
| © ChatGPT∞                         | NA   | 33   | 17     | 19             | 13       | 1         | 11        | 4        | 2         | 57.6% | 64.7%        |
| ♂ DeepSeeker-Inst                  | 33b  | 33   | 17     | 18             | 14       | 1         | 8         | 9        | 0         | 54.5% | 47.1%        |
| Claude-3-haiku                     | NA   | 28   | 22     | 16             | 10       | 2         | 11        | 11       | 0         | 57.1% | 50.0%        |
| 𝔍 DeepSeeker-1.5-Inst <sup>™</sup> | 7b   | 27   | 23     | 18             | 8        | 1         | 9         | 11       | 3         | 66.7% | <b>39.1%</b> |
| 🖸 Gemini 🖻                         | NA   | 19   | 31     | 13             | 6        | 0         | 10        | 18       | 3         | 68.4% | 32.3%        |

training or inference paradigms for LLMs for code focus on obtaining high quality datasets boosted with instruction-tuning, our result shows that existing LLMs still struggle with the concept of problem composition to tackle more complex problems. We hope future research can design novel training methods to tackle this limitation.

### 5.3 Problem Decomposition

Given our analysis and benchmark on combining different problems together, a nature follow-up would be to look at *problem decomposition* – decomposing larger problems into multiple subproblems. We start by selecting 50 HUMANEVAL problems and then follow our approach in Section 2 to decompose each original problem into two smaller subproblems, creating 100 problems in our DECOMPOSE benchmark.

Table 4 shows the results of selected LLMs on DECOMPOSE (the same set of LLMs as COM-BINE). We first observe that similar to the composition percentage in the COMBINE and COMBINE-NAIVE problems, LLMs do not achieve a high decomposition percentage. One possible interpretation is that current LLMs are trained to memorize or recover seen outputs in their training data, and when used for program synthesis, they cannot generalize the concepts from training data. This is demonstrated by not being able to solve smaller subproblems obtained from solved more difficult parent problems. On the other hand, we show that LLMs can sometimes solve both smaller subproblems even when the original parent problem is not solved (i.e., recomposition percentage). DECOMPOSE is akin to breaking the harder problem down into easier subproblems, which is related to planning in prior work [22]. We hope future work can again build on these insights to achieve the best of both worlds in being able to succesfully generalize difficult concepts into subproblems and adopting decomposing/planning to solve additional challenging problems.

#### 5.4 Tool Using

We further analyze the TOOL\_USE dataset, which contains pre-defined helper or auxiliary functions in addition to the main synthesis problem. Additionally, we construct TOOL\_USE-MAIN\_ONLY dataset, which contains the same set of problem as TOOL\_USE, except that the input to the LLM consists only of the main problem description without including any helpers. Using both datasets together, we can evaluate the ability of LLMs to use helper functions to solve more complex problem. We observe that compared to scenarios without any helper functions (average pass@1 of 28.6%), LLMs on average improve by 81.3% when provided with the helper functions. This is to be expected as the helper functions provides additional utilities in aiding to solve the more complex problem. However, this improvement is not uniform, as we see that the average improvement when given the auxiliary functions



Figure 9: pass@1 improvement from TOOL\_USE-MAIN\_ONLY to TOOL\_USE on selected instruction-following models and their base variants.

for instruction-following models is only 60.4% compared to the non-instruction-following LLMs' improvement of 122.0%.

Figure 9 show the detailed comparison between 10 instruction-following and their base LLMs on both the TOOL\_USE-MAIN\_ONLY and TOOL\_USE dataset. We observe that without the helpers, the instruction-following models significantly outperform their base LLMs. However, once the helpers are provided, this gap is drastically decreased, with cases even where the base models outperform their instruction-following counterparts. As real-world coding involves understanding, using, and then reusing existing functions across different places in the repository, being able to successfully leverage auxiliary methods is key. Current instruction-following LLMs are generally fine-tuned with data consisting of self-contained code snippets without the interaction and learning of function usages. This is further exacerbated by prior benchmarks, which mostly use self-contained functions, thus cannot expose the insufficient tool-using capability of such models. In EVOEVAL, with TOOL\_USE and TOOL\_USE-MAIN\_ONLY, we demonstrate this gap in evaluation and hope to inspire future research on this important aspect of code LLMs.

### 6 Related Work

Large language models for code. Starting with the general development of LLMs for general purpose tasks, developers have applied LLMs to perform code-related tasks by further training LLMs using collected code snippets from open-source repositories. Such LLMs include CODEX [10], PolyCoder [52], CodeT5 [48], CodeGen [34], InCoder [14], CodeLlama [40], StarCoder [26], StarCoder2 [29], DeepSeeker [16], etc. These LLMs can autoregressive complete code given the relevant prefix (e.g., docstrings for function completion). More recently, following the advancement in NLP, researchers have applied instruction-tuning methods to train code-specific LLMs that are well-versed in following instructions. Examples of such LLMs include CodeLlama-Inst [40] and DeepSeeker-Inst [16]. WizardCoder [30] instruction-tunes the model using Evol-Instruct to create more complex instructions. Magicoder [50] develops OSS-Instruct by synthesizing high quality instruction data from open-source code snippets. OpenCodeInterpreter [55] additionally leverages execution feedback for instruction-tuning in order to better support multi-turn code generation and refinement.

**Program synthesis benchmarking.** HUMANEVAL [10] and MBPP [4] are two of the most widely-used handcrafted code generation benchmarks complete with test cases to check for the correctness of LLM outputs. Building on these popular benchmarks, additional variants have been crafted including: HUMANEVAL+ [28] which improves the two benchmarks with more complete testcases; HUMANEVAL-X [54] which extends HUMANEVAL to C++, Javascript and Go; MultiPL-E [9] which further extends both HUMANEVAL and MBPP to 18 coding languages. Similarly, other benchmarks have been developed for specific domains: DS-1000 [25] and Arcade [53] for data science APIs; ODEX [49] for open-domain code generation covering a diverse range of libraries; CodeContests [27], APPS [17] and Live-CodeBench [20] for programming contests; ClassEval [13] for class-level generations, and

SWE-Bench [23] for real-world software engineering tasks. Different from prior benchmarks which require handcraft problems from scratch – high manual effort or scrape open-source repositories or coding contest websites – leading to unavoidable data leakage, EVOEVAL directly uses LLMs to *evolve* existing benchmark problems to create new complex evaluation problems. Furthermore, contrasting with the narrow scope of prior benchmarks (often focusing on a single type or problem, i.e., coding contests), EVOEVAL utilizes targeted transformation to evolve problems into different domains, allowing for a more holistic evaluation of program synthesis using LLMs.

### 7 Conclusion

We present EVOEVAL- a set of program synthesis benchmarks created by *evolving* existing problems into different target domains. We build on top of the popular HUMANEVAL benchmark to produce 828 problems across 7 different benchmarks for a holistic and comprehensive evaluation of LLM program synthesis ability. Our results on **51** LLMs show that compare to high performance on standard benchmarks, there is drastic drop in performance (on average 39.4%) when evaluated on EVOEVAL. Additionally, we observe significant ranking differences compared to previous leaderboards, indicating potential overfitting of popular LLMs on existing benchmarks. Throughout the paper, we provide additional insights, including the brittleness of instruction-following LLMs as well as problem composition and decomposition abilities. We hope EVOEVAL not only provides a valuable benchmarking suite for program synthesis but also inspires future code LLM builders to recognize the shown limitations of existing code LLMs and develop novel and targeted training approaches for code. We have open-sourced the EVOEVAL benchmarks, tools, and complete LLM generations available at https://github.com/evo-eval/evoeval

## 8 Acknowledgment

We thank Owen Colegrove for his help on starting this project and providing valuable feedback throughout, Jiawei Liu for providing helpful discussions and Yifeng Ding for his help in running experiments.

### References

- [1] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:*2305.10403, 2023.
- [2] Anthropic. Introducing claude 2.1. https://www.anthropic.com/news/claude-2-1/, 2023.
- [3] Anthropic. Introducing the next generation of claude. https://www.anthropic.com/ news/claude-3-family/, 2024.
- [4] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [6] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Xiaodong Deng Kai Dang, Yang Fan, Wenbin Ge, Fei Huang, Binyuan Hui, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Tianyu Liu, Keming Lu, Jianxin Ma, Rui Men, Na Ni, Xingzhang Ren, Xuancheng Ren, Zhou San, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Jin Xu, An Yang, Jian Yang, Kexin Yang, Shusheng Yang, Yang Yao, Jianwei Zhang Bowen Yu, Yichang Zhang, Zhenru Zhang, Bo Zheng, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Introducing qwen1.5. https://qwenlm.github.io/blog/qwen1.5/, 2023.

- [7] Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondřej Dušek. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source llms. arXiv preprint arXiv:2402.03927, 2024.
- [8] BudEcosystem. Code millenials 34b. URL [https://huggingface.co/ budecosystem/code-millenials-34b] (https://huggingface.co/budecosystem/ code-millenials-34b).
- [9] 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. Multipl-e: A scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions on Software Engineering*, 2023.
- [10] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374, 2021.
- [11] Yihe Deng, Weitong Zhang, Zixiang Chen, and Quanquan Gu. Rephrase and respond: Let large language models ask better questions for themselves, 2023.
- [12] Yinlin Deng, Chunqiu Steven Xia, Haoran Peng, Chenyuan Yang, and Lingming Zhang. Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models. In 32nd International Symposium on Software Testing and Analysis (ISSTA), 2023.
- [13] Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation. arXiv preprint arXiv:2308.01861, 2023.
- [14] Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Scott Yih, Luke Zettlemoyer, and Mike Lewis. Incoder: A generative model for code infilling and synthesis. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=hQwb-lbM6EL.
- [15] Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. Program synthesis. Foundations and Trends® in Programming Languages, 4(1-2):1–119, 2017. ISSN 2325-1107. doi: 10. 1561/2500000010. URL http://dx.doi.org/10.1561/2500000010.
- [16] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming–the rise of code intelligence. *arXiv preprint arXiv:*2401.14196, 2024.
- [17] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps. *NeurIPS*, 2021.
- [18] Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15, 2002.
- [19] HuggingFace. Hugging face, 2022. https://huggingface.co.
- [20] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. arXiv preprint, 2024.
- [21] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [22] Xue Jiang, Yihong Dong, Lecheng Wang, Qiwei Shang, and Ge Li. Self-planning code generation with large language model. *arXiv preprint arXiv:2303.06689*, 2023.

- [23] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
- [24] Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. Measuring compositional generalization: A comprehensive method on realistic data. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum? id=SygcCnNKwr.
- [25] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pp. 18319–18345. PMLR, 2023.
- [26] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you!, 2023.
- [27] 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. Competition-level code generation with alphacode. *Science*, 2022. URL https://www.science.org/doi/abs/10.1126/science.abq1158.
- [28] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=1qvx610Cu7.
- [29] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.
- [30] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- [31] William M McKeeman. Differential testing for software. *Digital Technical Journal*, 10(1): 100–107, 1998.
- [32] Microsoft Research. Phi-2: The surprising power of small language models. https://www.microsoft.com/en-us/research/blog/ phi-2-the-surprising-power-of-small-language-models/, 2023.
- [33] Mistral AI team. Mixtral of experts a high quality sparse mixture-of-experts. https: //mistral.ai/news/mixtral-of-experts/, 2023.

- [34] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=iaYcJKpY2B\_.
- [35] OpenAI. Chatgpt: Optimizing language models for dialogue. https://openai.com/ blog/chatgpt/, 2022.
- [36] OpenAI. Gpt-4 technical report. ArXiv, abs/2303.08774, 2023.
- [37] phind team. Beating gpt-4 on humaneval with a fine-tuned codellama-34b. https: //www.phind.com/blog/code-llama-beats-gpt4, 2023.
- [38] Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, and Nathan Cooper. Stable code 3b. URL [https://huggingface.co/stabilityai/ stable-code-3b] (https://huggingface.co/stabilityai/stable-code-3b).
- [39] Martin Riddell, Ansong Ni, and Arman Cohan. Quantifying contamination in evaluating code generation capabilities of language models. *arXiv preprint arXiv:2403.04811*, 2024.
- [40] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950, 2023.
- [41] Kensen Shi, Joey Hong, Yinlin Deng, Pengcheng Yin, Manzil Zaheer, and Charles Sutton. Exedec: Execution decomposition for compositional generalization in neural program synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=oTRwljRgiv.
- [42] Jiangwen Su. Code millenials 34b. URL [https://huggingface.co/ uukuguy/speechless-codellama-34b-v2.0] (https://huggingface.co/uukuguy/ speechless-codellama-34b-v2.0).
- [43] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [44] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv*:2403.08295, 2024.
- [45] Xwin-LM Team. Xwin-lm. https://github.com/Xwin-LM/Xwin-LM, 2023.
- [46] Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv*:2309.11235, 2023.
- [47] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- [48] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8696–8708, 2021.
- [49] Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. Execution-based evaluation for open-domain code generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 1271–1290, 2023.
- [50] Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*, 2023.

- [51] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. Automated program repair in the era of large pre-trained language models. In *Proceedings of the 45th International Conference on Software Engineering (ICSE 2023). Association for Computing Machinery*, 2023.
- [52] Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of large language models of code. In *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pp. 1–10, 2022.
- [53] Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Alex Polozov, and Charles Sutton. Natural language to code generation in interactive data science notebooks. 2022.
- [54] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*, 2023.
- [55] Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhu Chen, and Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv preprint arXiv:2402.14658*, 2024.

### **A** Evaluation LLMs

#### A.1 Evaluated LLMs

Table 5 shows the overview of the 51 LLMs we evaluated in our work. For any LLMs which provide their open-source weights, we directly obtain them from huggingface model hub [19]<sup>3</sup>. For any close-sourced LLMs, we directly access their model endpoints using their providers. For more detail on the access of each LLM, please check our repository: https://github.com/evo-eval/evoeval

### A.2 Detailed Evaluation Setup

```
Please complete the following code snippet.
def transform_canvas(canvas: str) -> str:
    11 11 11
    You have an canvas containing either '#' (representing a wall), '-' (
    \hookrightarrow representing
    an empty space), or 'P' (representing the point at which a painter starts).
    \hookrightarrow The painter
    can move horizontally on the canvas and paints all empty spaces he encounters
    with '*' without crossing or hitting the walls.
    The task is to return an updated canvas with all the accessible spaces
    \hookrightarrow painted,
    keeping wall configuration and unaccessible spaces same. If the canvas
    \hookrightarrow contains no painter 'P',
    return the canvas as it is. If there are more than one 'P' or the number of
    \hookrightarrow painted space divides the empty spaces evenly, return 'Invalid canvas'.
    Examples:
    >>> transform_canvas('P----#-----#-----')
    'P***#-----#-----#-----'
    >>> transform_canvas('--#-P#----#---#---*')
    'Invalid canvas'
    >>> transform_canvas('----#--P--#-----')
    '----#**P**#----#----'
    >>> transform_canvas('----#--P---#P----')
    'Invalid canvas'
```

Figure 10: Example input prompt for GPT-4

**LLM generation.** As mentioned in Section 4, we report the pass@1 score for each LLM on our dataset generated using greedy decoding (i.e., sampling with temperature = 0). For each LLM, we provide a specific input prompt depending on the model type. For base LLMs (i.e., not instruction-following variants), we use only the function headers as input. For instruction-following, we make the best effort to follow examples provided by each model maker on the exact instruction and format to use at the time of writing. Specifically, for instruction-following LLMs, we ask the model to return the code snippet wrapped by code blocks (i.e., ```). Figure 10 shows an example input for GPT-4 on a CREATIVE problem.

Furthermore, we also provide a custom sanitization script adopted from EVALPLUS [28] which parses the raw LLM outputs for code block parsing (e.g., removing ``` indicators for

<sup>&</sup>lt;sup>3</sup>For certain LLMs, we may use the vLLM inference library for more efficient generation

Table 5: Detailed overview of evaluated models. Model ID indicates either the API endpoint name or huggingface model name used for the particular model. Available Weights indicate whether the model is evaluated by accessing a close-sourced API endpoint or ran locally with provided weights. Note: Detailed evaluated by accessing a close-sourced API endpoint or ran locally with provided weights. Note: Detailed evaluated by accessing a close-sourced API endpoint or ran locally with provided weights. Note: Detailed evaluated by accessing a close-sourced API endpoint or ran locally with provided weights. Note: Detailed evaluated by accessing a close-sourced evaluated by accessing a close-sourced evaluated by accessing evaluated by accessing a close-sourced evaluated by accessing evaluated by accessing evaluated by accessing evaluated by accessing evaluated evaluated by accessing evaluated evaluated evaluated by accessing evaluated evaluat

|  | Size                    | Model ID  | Available Weights                |
|--|-------------------------|---|----------------------------------|
| @ GPT-4-Turbo=[36]                     | NA                      | gpt-4-0125-preview  | X                                |
| \$ GPT-4 <sup>[2]</sup> [36]           | NA                      | gpt-4-0613  | ×                                |
| <sup>®</sup> ChatGPT <sup>™</sup> [35] | NA                      | gpt-3.5-turbo-0125  | ×                                |
| Claude-3 <sup>[3]</sup>                | NA                      | claude-3-opus-20240229  | ×                                |
| Claude-3-haiku <sup>[3]</sup>          | NA                      | claude-3-haiku-20240307   | ×                                |
| Claude-2 <sup>[2]</sup>                | NA                      | claude-2.1  | ×                                |
| • Gemini <sup>®</sup> [43]             | NA                      | gemini-1.0-pro  | ×                                |
| • PaLM-2 <sup>[1]</sup>                | NA                      | text-bison-001  | ×                                |
| ♂ DeepSeeker-Inst∞[16]                 | 33b<br>6.7b<br>1.3b     | deepseek-ai/deepseek-coder-33b-instruct<br>deepseek-ai/deepseek-coder-6.7b-instruct<br>deepseek-ai/deepseek-coder-1.3b-instruct                         | J<br>J<br>J                      |
| ♂ DeepSeeker [16]                      | 33b<br>6.7b<br>1.3b     | deepseek-ai/deepseek-coder-33b-base<br>deepseek-ai/deepseek-coder-6.7b-base<br>deepseek-ai/deepseek-coder-1.3b-base                                     | \$<br>\$<br>\$                   |
| ▼ DeepSeeker-1.5-Inst.@[16]            | 7b                      | deepseek-ai/deepseek-coder-7b-instruct-v1.5   | 1                                |
| VeepSeeker-1.5 [16]                    | 7b                      | deepseek-ai/deepseek-coder-7b-base-v1.5   | 1                                |
| ∞ CodeLlama-Inst@[40]                  | 70b<br>34b<br>13b<br>7b | codellama/CodeLlama-70b-Instruct-hf<br>codellama/CodeLlama-34b-Instruct-hf<br>codellama/CodeLlama-13b-Instruct-hf<br>codellama/CodeLlama-7b-Instruct-hf | <i>v</i><br><i>v</i><br><i>v</i> |
| ∞ CodeLlama [40]                       | 70b<br>34b<br>13b<br>7b | codellama/CodeLlama-70b-Python-hf<br>codellama/CodeLlama-34b-Python-hf<br>codellama/CodeLlama-13b-Python-hf<br>codellama/CodeLlama-7b-Python-hf         | \<br>\<br>\<br>\                 |
| WizardCoder [30]                       | 34b                     | WizardLM/WizardCoder-Python-34B-V1.0  |                                  |
| WizardCoder-1.1 <sup>[9]</sup> [30]    | 33b                     | WizardLM/WizardCoder-33B-V1.1   | ✓                                |
| XwinCoder [45]                         | 34b                     | Xwin-LM/XwinCoder-34B   | 1                                |
| Phind-CodeLlama-2 [37]                 | 34b                     | Phind/Phind-CodeLlama-34B-v2  | 1                                |
| Code Millenials <sup>[8]</sup>         | 34b                     | budecosystem/code-millenials-34b  | 1                                |
| Speechless-CL <sup>[]</sup> [42]       | 34b                     | uukuguy/speechless-codellama-34b-v2.0   | 1                                |
| Magicoder-s-DS <sup>[50]</sup>         | 6.7b                    | ise-uiuc/Magicoder-S-DS-6.7B  | 1                                |
| Magicoder-s-CL=[50]                    | 7b                      | ise-uiuc/Magicoder-S-CL-7B  | 1                                |
| StarCoder2 [29]                        | 15b<br>7b<br>3b         | bigcode/starcoder2-15b<br>bigcode/starcoder2-7b<br>bigcode/starcoder2-3b  | \<br>\<br>\                      |
| StarCoder [26]                         | 15b                     | bigcode/starcoder   | 1                                |
| Mixtral-Inst=[33]                      | 8x7b                    | mistralai/Mixtral-8x7B-Instruct-v0.1  | 1                                |
| Mistral-Inst-v02 <sup>[21]</sup>       | 7b                      | mistralai/Mistral-7B-Instruct-v0.2  | 1                                |
| Mistral-Inst🕬[21]                      | 7b                      | mistralai/Mistral-7B-Instruct-v0.1  | 1                                |
| Mistral [21]                           | 7b                      | mistralai/Mistral-7B-v0.1   | 1                                |
| OpenChat <sup>[9]</sup> [46]           | 7b                      | openchat/openchat-3.5-0106  | 1                                |
| stable-code [38]                       | 3b                      | stabilityai/stable-code-3b  | 1                                |
| O Gemma-Inst.∞ [44]                    | 7b                      | google/gemma-7b-it  | 1                                |
| <b>O</b> Gemma [44]                    | 7b<br>2b                | google/gemma-7b<br>google/gemma-2b  | <i>J</i><br><i>J</i>             |
| Phi-2 [32]                             | 2.7b                    | microsoft/phi-2   | 1                                |
| Qwen-1.5 <sup>[6]</sup>                | 72b<br>14b<br>7b        | Qwen/Qwen1.5-72B-Chat<br>Qwen/Qwen1.5-14B-Chat<br>Qwen/Qwen1.5-7B-Chat  | <i>s</i><br><i>s</i>             |
| Qwen <sup>[5]</sup>                    | 14b<br>7b               | Qwen/Qwen-14B-Chat<br>Qwen/Qwen-7B-Chat   | 1<br>1                           |

instruction-following models) and end-of-string identifiers (e.g., removing tokens like </s>). Each model generated output is passed into the sanitization script and the evaluation occurs on the sanitized outputs.

**Oracle.** To evaluate the functional correctness of each LLM synthesized solution, we use differential testing by comparing the model output with the groundtruth output on a set of testcase inputs. We build our evaluation framework on top of the EVALPLUS evaluation script used for HUMANEVAL and HUMANEVAL+ benchmark which evaluates multiple problems and solutions in parallel for efficiency. For each testcase, we perform exact matching or check if the output is within an absolute difference threshold of  $10^{-6}$  if the output is a floating point type. We additionally implement our evaluation script by recursively checking the type and performing the appropriate comparison (e.g., dictionary outputs are first length checked for equivalence and then matching is done for each value and key). Furthermore, we also implement custom oracles for specific problems where there could be multiple solutions or simple tolerance or exact matching cannot fully guarantee correctness. We refer the reader again to our repository https://github.com/evo-eval/evoeval which contains the full implementation of each of our custom oracle. Additionally, we also use timeout as another evaluation method. Our setting again follows EVALPLUS default setup where the timeout per problem is defined as  $T = max(T_{max}, f \times t_{gt})$  with default values of  $T_{max} = 1000ms$ , f = 4 and  $t_{gt}$  defined as the measured groundtruth solution time to produce the correct output. All timeout related factors can be adjusted to account for variance on different underlying machine and hardware.

## **B** Transformation Prompts

Figure 11: Prompt for DIFFICULT

Figure 12: Prompt for CREATIVE

Here we provide the exact targeted transformation prompts used to evolve existing benchmark problems into each of our transformation prompts. Figure 11, 12, 13, 14, 15, 16, 17 and 18 shows the prompt for DIFFICULT, CREATIVE, SUBTLE, COMBINE, TOOL\_USE, VERBOSE, CONCISE and DECOMPOSE respectively.

Figure 13: Prompt for SUBTLE

Figure 14: Prompt for COMBINE

Figure 19 and 20 show the refinement and I/O extraction/fixing prompt used in EVOEVAL. The refinement prompt is used to refine the origin generated problem when inconsistency is detected (see Section 2). The extraction prompt is used to initially obtain a set of testcases from the problem docstring used for self-consistency evaluation. We further use an I/O fixing prompt (also in Figure 20) to fix any examples in the docstring which do not contain the right output (as computed by the groundtruth generated by GPT-4).

### **C** Example Problems in EVOEVAL

Here we demonstrate a few example problems across the benchmarks in EVOEVAL and corresponding GPT-4 solution which cannot solve the problem. Figure 21, 22, 23, 24, 25,

```
Here is an example coding problem:
{problem}
Please come up with a new problem which uses helper functions to solve the
    \hookrightarrow problem.
The new problem should contain the following:
first: one or more helper functions
second: the main problem description consist of the function header and docstring
The main problem description should not refer to the helper function(s) in any
    \hookrightarrow way.
The helper function(s) should implement simple parsing or checking logic.
To solve the main problem, one should also use additional complex logic than just
    \hookrightarrow calling the helper function(s).
Avoid problems on simple math concepts such as prime, palindrome, anagrams,
    \hookrightarrow factorial
Avoid concepts like emails, string or parsing-based problems
Please return the full implementation of the helper function(s) and the main
    \hookrightarrow problem description (not the implementation) in the same format as the
    \hookrightarrow example problem (
```

Figure 15: Prompt for TOOL\_USE

Figure 16: Prompt for VERBOSE

Below is a coding problem

{problem}

Figure 17: Prompt for VERBOSE

26 and 27 show such example for the EVOEVAL DIFFICULT, CREATIVE, SUBTLE, COMBINE, TOOL\_USE, VERBOSE and CONCISE respectively.

#### Figure 18: Prompt for DECOMPOSE

```
Below is a coding problem
{problem}
Ensure logical coherence in the given problem.
Improve the docstring's clarity and conciseness.
Fix missing or helpful imports.
Include example input/output if absent.
Return the modified problem in the same format as the example problem (i.e.,
```

Figure 19: Refinement prompt

```
Here is a function header with docstring:
{problem}
Please extract the example raw input argument and expected output from the
     \hookrightarrow docstring.
If there are no example input and output please provide new ones.
Return each pair of input and output as assertions in this format:
assert {function_name}({{the_first_input_example}} == {{the_first_output_example}
     \hookrightarrow }}
assert {function_name}({{the_second_input_example}} == {{
    \hookrightarrow the_second_output_example}}
. . .
Here is a problem with docstring:
{problem}
Some example inputs and outputs in the docstring may be wrong. Please correct
    \hookrightarrow them according to the provided correct assertions below, and ensure that
    \hookrightarrow the correct example inputs and outputs <code>assertionsReturn</code> the <code>revised</code>
    problem in the same format as the original problem (i.e.,
```

Figure 20: Input extraction and fixing prompts

```
def common(l1: list, l2: list, n: int):
    .....
   Return the n longest strings, sorted by increasing length that are common in
   \hookrightarrow two lists.
   However, in the case of a tie, prioritize the string that appears first in
    \hookrightarrow list1.
   >>> common(["apple", "banana", "cherry", "date", "elderberry"], ["banana", "
    \hookrightarrow date", "fig", "grape", "elderberry"], 3)
   ["banana", "date", "elderberry"]
>>> common(["cat", "dog", "elephant", "fox", "goat"], ["dog", "fox", "cat", "
    \hookrightarrow horse", "iguana", "goat"], 3)
    ["cat", "dog", "goat"]
   ["absolutely"]
    ......
   common_strings = [string for string in 11 if string in 12]
    common_strings.sort(key=len) # does not adjust for tie sorting scenario.
   return common_strings[-n:]
```

Figure 21: GPT-4 failing solution on problem EVOEVAL/58 in DIFFICULT

```
def crossword_validator(grid: list, words: list):
    0.0.0
    Consider a 2D grid representing a square crossword puzzle. Each cell of the
    \hookrightarrow grid contains a single character.
    Determine if all the words from a provided list of words can be found in the
    \hookrightarrow grid.
    Words can be found left-to-right (in a row), right-to-left (in a row),
    top-to-bottom down (in a column), and bottom-to-top (in a column).
    We assume that words are only found in contiguous cells and words do not '
    \hookrightarrow wrap' around the grid.
    Note: the grid is a square (N x N) and all words are lowercase.
    Parameters:
    grid: list of list of char - where each char is a lowercase letter.
    words: list of str - where each str is a lowercase word.
    Returns:
    boolean - return True all words are found in the grid, False otherwise.
    >>> crossword_validator([['c', 'a', 't'], ['d', 'o', 'g'], ['d', 'o', 'p']],
    \hookrightarrow ['cat', 'dog', 'pop'])
    False
    >>> crossword_validator([['c', 'a', 't'], ['d', 'o', 'g'], ['d', 'o', 'p']],
    \hookrightarrow ['cdd', 'dog', 'pod'])
    True
    .....
    n = len(grid)
    for word in words:
        found = False
        for i in range(n):
            # does not correctly check if a word can be found as a sublist of
    \hookrightarrow column/grid
            if ''.join(grid[i]) == word or ''.join(grid[i][::-1]) == word:
                 found = True
                 break
             if ''.join([grid[j][i] for j in range(n)]) == word or ''.join([grid[j
    \hookrightarrow ][i] for j in range(n)][::-1]) == word:
                 found = True
                 break
        if not found:
            return False
    return True
```

Figure 22: GPT-4 failing solution on problem EVOEVAL/21 in CREATIVE

```
from typing import List
def below_one(operations: List[int]) -> bool:
    0.0.0
    You're given a list of deposit and withdrawal operations on a bank account
    \hookrightarrow that starts with
    one balance. Your task is to detect if at any point the balance of account
    \hookrightarrow falls below one, and
    at that point function should return True. Otherwise it should return False.
    >>> below_one([1, 2, 3])
    False
    >>> below_one([1, 2, -4, 5])
    True
    0.0.0
    # incorrect starts balance at 0 instead of 1
    balance = 0
    for operation in operations:
        balance += operation
        if balance < 1:</pre>
            return True
    return False
```

Figure 23: GPT-4 failing solution on problem EVOEVAL/3 in SUBTLE

```
def grade_with_same_chars(grades: list, s0: str, s1: str):
    A teacher has a special grading system where she assigns her students letter
    \hookrightarrow grades
    based on their GPA according to the following table:
                        GPA
                             Letter grade
              4.0
                                   A+
            > 3.7
                                   Α
            > 3.3
                                   A-
            > 3.0
                                   B+
            > 2.7
                                   R
            > 2.3
                                  B-
            > 2.0
                                  C+
            > 1.7
                                   С
            > 1.3
                                   C-
                                  D+
            > 1.0
            > 0.7
                                   D
            > 0.0
                                   D-
              0.0
                                   E
    However, this year she adds an extra condition: the students can only receive
    \hookrightarrow a grade
    if their name (s0) and the course title (s1) have the same characters. The
    \hookrightarrow function
    receives a list of GPAs and two strings (s0 and s1) and returns a list of
    \hookrightarrow grades
    for the students that pass the condition, or an empty list if no student
    \hookrightarrow passes.
    Example:
    grade_with_same_chars([4.0, 3, 1.7, 2, 3.5], 'abcd', 'dcba') ==> ['A+', 'B',
    \hookrightarrow 'C-', 'C', 'A-']
    grade_with_same_chars([4.0, 3, 1.7, 2, 3.5], 'aaabcd', 'dddcba') ==> ['A+', '
    \hookrightarrow B', 'C-', 'C', 'A-']
    grade_with_same_chars([4.0, 3, 1.7, 2, 3.5], 'abcd', 'dcbae') ==> []
    # does not correctly check characters are same
    if sorted(s0) != sorted(s1):
        return []
    grading_system = {4.0: 'A+', 3.7: 'A', 3.3: 'A-', 3.0: 'B+', 2.7: 'B', 2.3: '
    ↔ B-', 2.0: 'C+', 1.7: 'C', 1.3: 'C-', 1.0: 'D+', 0.7: 'D', 0.0: 'D-', 0.0:
    \hookrightarrow 'E'}
    grades_assigned = []
    for gpa in grades:
        for key in sorted(grading_system.keys(), reverse=True):
             if gpa >= key:
                 grades_assigned.append(grading_system[key])
                 break
    return grades_assigned
```

Figure 24: GPT-4 failing solution on problem EVOEVAL/28 in COMBINE

```
def is_valid_passport(passport: str) -> bool:
    ......
    Checks whether a given string is a valid passport number.
    A valid passport number is defined as having exactly 9 alphanumeric
    \hookrightarrow \text{ characters.}
    >>> is_valid_passport('P12345678')
    True
    >>> is_valid_passport('P123456789')
    False
    >>> is_valid_passport('P12345G78')
    False
    .....
    if len(passport) != 9:
        return False
    if not passport.isalnum():
        return False
    return True
from typing import List, Tuple
def process_passport_data(data: str) -> List[Tuple[str, str]]:
    Parse a string containing passport data and return a list of tuples.
    Each tuple contains the passport number and the name of the passport holder.
    The passport data string is in the format of 'PassportNumber:Name'
    separated by lines.
    Only process valid passport data. A valid passport number is defined as
    \hookrightarrow having exactly 9 alphanumeric characters.
    >>> process_passport_data('P12345678:John Doe\nX12345678:Jane Doe')
    [('P12345678', 'John Doe')]
    >>> process_passport_data('P12345678:John Doe\nA123456789:Jane Doe')
    [('P12345678', 'John Doe')]
    0.0.0
    valid_passports = []
    lines = data.split('\n')
    for line in lines:
        (passport, name) = line.split(':') # directly splits without checking for
        validity
    \hookrightarrow
        if is_valid_passport(passport):
            valid_passports.append((passport, name))
    return valid_passports
```

Figure 25: GPT-4 failing solution on problem EVOEVAL/23 in TOOL\_USE

```
def starts_one_ends(n):
    .....
    This function, given a user-provided positive integer 'n', calculates and
    \hookrightarrow returns the total count of n-digit positive integers where 1 appears
    \hookrightarrow either at the start or the end of the number.
    For instance, if the input is 2, the function will count all the two-digit
    \hookrightarrow numbers that either start or end with 1, like 10, 11, 21, 31, ..., 91,
    \hookrightarrow and 12, 13, ..., 19.
    It is important to note that this function only accepts positive integers. If
    \hookrightarrow a non-integer or non-positive integer is provided, the function may not
    \hookrightarrow behave as expected.
    Args:
    n : int
        A positive integer that represents the number of digits in the numbers to
    \hookrightarrow be counted.
    Returns:
    int
        The count of the numbers of n-digit positive integers that start or end
    \hookrightarrow with 1.
    Examples:
    The function can be called in the following way:
        starts_one_ends(2)
    This example call would return the count of two-digit numbers where 1 appears
    \hookrightarrow at the start or at the end.
    .....
    if n == 1:
        return 1
    else:
        # incorrectly adds 2 to the final answer
        return 18 * 10 ** (n - 2) + 2
```

Figure 26: GPT-4 failing solution on problem HUMANEVAL/83 in VERBOSE

```
def circular_shift(x, shift):
    """Performs a circular right shift on the digits of x as per shift value and
    ↔ returns the result in string format.
    Returns reversed digits if shift exceeds number of digits.
    >>> circular_shift(12, 1)
    "21"
    >>> circular_shift(12, 2)
    "12"
    """
    x = str(x)
    # does not check shift exceeds string length and return reversed digits
    shift = shift return x[-shift:] + x[:-shift]
```

