# SELFEVOLVE: A Code Evolution Framework via Large Language Models

**Shuyang Jiang**[1], **Yuhao Wang**[1], **Yu Wang**[1,2*]

[1]Shanghai Jiao Tong University
[2]Shanghai AI Laboratory
{jiangshuyang,colane,yuwangsjtu}@sjtu.edu.cn

## Abstract

Large language models (LLMs) have already revolutionized code generation, after being pretrained on publicly available code data. However, while various methods have been proposed to augment LLMs with retrieved knowledge and enhance the quality of code generation, the performance of these retrieval-based methods is limited by the strength of the retrievers used. In addition, while LLMs show great emergent ability, they still struggle to produce the correct code in one turn. To address these challenges, we propose a novel two-step pipeline, called SELF-EVOLVE, that leverages LLMs as both knowledge providers and self-reflective programmers. Unlike retrieval-based methods, SELFEVOLVE obtains the knowledge from input prompts and generates intermediate code based on the generated knowledge. After that, SELFEVOLVE asks LLM to act as an expert programmer to perform debugging for the generated code. This is achieved by receiving the error message from the interpreter, without requiring special test cases for correctness verification. We evaluate SELFEVOLVE on three code generation datasets, including DS-1000 for data science code, HumanEval for software engineering code, and TransCoder for C++-to-Python translation. Our empirical experiments show that SELFEVOLVE outperforms strong baselines by a significant margin on all datasets. We also conduct exhaustive analytical experiments to validate the effectiveness of the two stages of SELFEVOLVE, and find that both are superior to other prompting-based methods. Further scalability analysis demonstrates that SELFEVOLVE can be adapted to other more advanced models, such as GPT-4, and bring consistent efficacy improvement.

## 1 Introduction

Code generation functions as a crucial and challenging component of various applications [2, 19, 20, 26]. However, the performance of large language models (LLM) on diverse tasks and domains has substantially improved as the pretraining corpus expands. As a result, LLM has become the preferred model for code generation [9, 28]. In fact, LLM performs much better than previous deep neural models dedicated to generating code [12, 26, 32]. Meanwhile, previous methods have been augmented by LLM's ability to digest various prompt contents and perform text generation [3, 22, 31]. Various auxiliary augmentation signals have been added to the prompt to obtain more accurate code [8, 18, 44, 56]. However, most prior work usually obtains such signals via an external retriever and a large knowledge base. They leverage the problem description or natural language intents to retrieve relevant knowledge, including similar code snippets [35, 36], API documentation [49, 56], or focal methods [25]. Despite their success, retriever models can suffer from domain mismatch when adapting to different tasks, requiring finetuning or even training from scratch on the target domain,

---

*Corresponding author

which limits their generality. Moreover, current retrievers are not well-suited for semi-structured knowledge items like library documentation, which can result in poor retrieved results.

To avoid domain mismatch and inaccurate retrieval results, we propose a two-stage paradigm, SELFEVOLVE, which treats LLM itself as a knowledge source. Previous work has demonstrated that LLM has encoded various domain knowledge [14] and can be treated as a large knowledge base [1, 39]. Therefore, SELFEVOLVE chooses to prompt LLM to generate multi-form necessary knowledge by itself. In particular, we prompt the language models to extract necessary knowledge from trial solutions or problem intents (§3.2), depending on whether the problem intents contain explicit demands. This process excludes the intervention of retrievers since generating concrete knowledge based on roughly encoded ones in LLM parameters is easier than searching them in a large database with vague natural language statements. To the best of our knowledge, SELFEVOLVE is the first LLM-driven self-augmented code generation framework. Furthermore, inspired by the fact that human programmers rely on both related knowledge and a debugger to ensure implementation correctness, we inject an automatic refinement mechanism. This refinement mechanism teaches language models to depend on an executor like a Python interpreter to correct the preliminary code. We construct a runnable program from generated code and test cases extracted from the problem description (§3.2) and execute it to obtain either pass or error messages, which serve as correction feedback. Compared to prompting LLM to generate test cases like CodeT [8], which may produce incorrect samples, SELFEVOLVE maintains the correctness of test cases. Additionally, we do not grab evaluation samples from the test set like the recently proposed Self-Debugging [10], which hardly generalizes to daily coding scenarios. Instead, the example cases in the problem description appear in coding tasks mostly and describe the behavior that the programmer needs to accomplish with little ambiguity. Leveraging these authentic and common test cases makes SELFEVOLVE a reliable and general method for self-augmented code generation.

We primarily build SELFEVOLVE using `gpt-3.5-turbo` [1] (ChatGPT) and evaluate its performance on various tasks, including the data science code generation task DS-1000 [26], the general code generation task HumanEval [9], and the C++-to-Python translation task TransCoder [43]. Extensive experiments show that SELFEVOLVE achieves a significant improvement in execution-based measurements compared to strong DocPrompting [56] and Self-Debugging [10] baselines on the data science code generation task (§4.2). On HumanEval, SELFEVOLVE still outperforms each strong baseline, and the self-refinement module brings noticeable performance improvements for the base LLM after using self-generated knowledge (§4.2). Even on the code translation, which is a much simpler task for ChatGPT, SELFEVOLVE still brings considerable improvement (§4.2). Furthermore, our analysis studies indicate that SELFEVOLVE can provide more accurate knowledge than retrieval-based methods (§4.3), generalize to various datasets with only a few debugging turns (§4.3), and scales easily to more powerful models like GPT-4 (§4.3). Finally, we use two intuitive cases to demonstrate how the two stages of SELFEVOLVEimprove the generated code. These cases illustrate the effectiveness of the model in generating high-quality code and highlight its potential for a range of applications.

## 2   Related Work

**Augmented code generation**   In addition to problem descriptions and code snippets, many works provide auxiliary information to generate code. Before the era of LLM, researchers trained an encoder-decoder model that aimed to generate code based on a programmatic environment and function documentation [21]. JuPyT5 [7] conditions on Jupyter notebooks' context cells to generate data science code. Recently, large language models (LLM), pretrained on a variety of corpus, have enabled an in-context learning pipeline for zero-shot or few-shot generation. Haluptzok et al. [17] introduced a method to solve programming puzzles with synthetic puzzles and solutions generated by LLM. Parvez et al. [35] augmented code generation models with retrieved similar code snippets. In contrast, API-specific documents retrieved via a CodeT5 [47] retriever serve as additional information in the prompt [56]. However, their methods involve retrieving semi-structured knowledge items, whose performance is bottlenecked by current retriever models. Moreover, it is hard for retrievers to adapt to the target domain when the corpus for finetuning the retriever is inaccessible. Compared to theirs, LLM is more suitable for bridging the gap between domains than any small retriever, which provides more accurate knowledge. Moreover, our method does not require domain-specific finetuning, offering higher accessibility and generality.

---
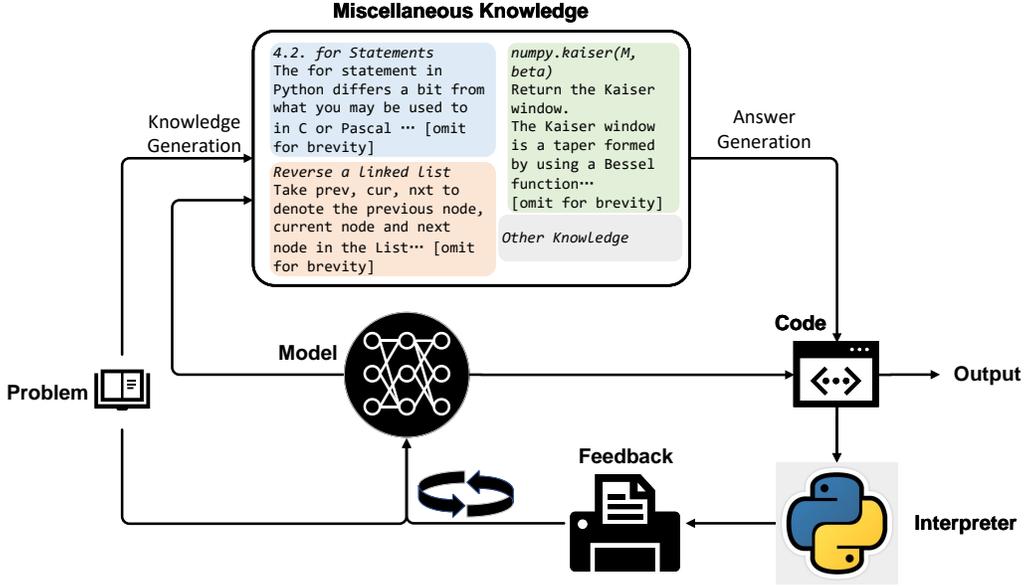
[1] https://platform.openai.com/docs/api-reference/chat

Figure 1: The SELFEVOLVE pipeline. LLMs first generate corresponding knowledge for the related problem, then generate the trial answer conditioned on the knowledge. The iterative refinement step uses test cases and generated code snippets to form executable programs and then prompts LLM to refine the answer code based on the feedback thrown by the interpreter.

**Automatic code refinement** Language models often output unreliable information with high confidence [24]. This unreliability is reflected in bug code snippets when generating code. Some previous works have trained exclusive models dedicated to repairing incorrect code, which only accept the bug code as input [16, 51]. Other works have fused auxiliary information obtained from execution, such as the stack trace [15, 45] and error information raised by a compiler [50]. In recent work, LLMs have been leveraged to act as "teachers" to fix bugs hidden in the code. CodeT [8] uses synthetic test cases obtained through Codex [9] to select correct codes. However, CodeT is applied in an unnatural scenario where all coding problems are formatted as writing a function with certain inputs. Madaan et al. [30] iteratively refines the output with model-generated feedback in multiple tasks, except for code generation. Chen et al. [10] uses similar methods in code generation, but their work peeps one test case from the ground truth, which is impossible when solving real problems. Without exposing test cases, our method is much more flexible and closer to a real coding scenario.

## 3   SELFEVOLVE: Code Evolution via Large Language Models

This section first briefly introduces the code generation paradigm inspired by natural programming practice. We then present the concept of SELFEVOLVE, a two-step method that utilizes language models as a self-bootstrapping knowledge enhancer and an expert programmer with self-reflection, without external models or knowledge bases. After presenting the concept, we delve into the two primary components of SELFEVOLVE, which are built up to form a fully LLM-driven generation pipeline without the need for fine-grained prompt designs or further finetuning steps. The overall pipeline of our method is presented in Figure 1.

### 3.1   BackGround

**Code generation formulation** Given a problem description written in natural language $d$, and code context $c$, the autoregressive language model $p_\theta$ predicts the solution as

$$P(Y) = \prod_{i=1..n} p_\theta(Y_i|Y_{<i}, X), Y_{<1} = \emptyset \tag{1}$$

where $n$ is the prediction length and $X = [d; c]$ is the concatenation of $d$ and $c$.

3

**Two-step code generation pipeline**   Conditioning solely on problem description for generation is still hard for LLM. Inspired by most programmers who often refer to knowledge documentation [42] and struggle to debug with current tools [34], we divide prompt-based code generation methods into two steps. The first step prompts language models to comprehend extra knowledge and task-specific instructions [4, 23, 31, 37], while the next one teaches models to revise the generated code solution through feedback from humans or an oracle instructor. In the two-step pipeline, the second generation step never deteriorates the intermediate output of the first step. Therefore, these two steps follow a topological order in terms of optimization, and can be optimized in order and fused together.

## 3.2   SELFEVOLVE: the Proposed Two-step Pipeline

Based on the above analysis, we propose SELFEVOLVE, which improves both steps by enabling generated code to evolve progressively using only a large language model, without requiring any learning. SELFEVOLVE generates code by conditioning on the knowledge in the prompt, as previous work has done. However, the knowledge is generated by LLM instead of being retrieved from external knowledge bases. After obtaining the output of the first step, SELFEVOLVE uses LLM to iteratively revise the generated code. This process follows Chen et al. [10] to correct code errors by leveraging feedback from a code executor, but does not necessitate the use of specific test cases.

**Generating knowledge with language models**   Conditioning a language model on the knowledge in the prompt is crucial, yet challenging. Given $m$ knowledge items $K[1..m]$, the language model predicts the next token to generate the final code solution:

$$P(Y|K) = \prod_{i=1..n} p_\theta(Y_i|Y_{<i}, X, K), Y_{<1} = \emptyset \qquad (2)$$

Knowledge can be retrieved via a sparse retriever [41] or a dense retriever [13, 40] as such:

$$K := \arg\max_{K \subset B} P(K|X, B) \qquad (3)$$

where $B$ is the whole database. However, the performance of current retriever models may be limited, resulting in $K$ containing irrelevant knowledge items that add noise to LLM and harm the generation results. A widely-used approach to mitigate this problem is to retrieve as much knowledge as possible [3] to cover the necessary items. However, this method places demands on the ability of LLM to process long texts, which is still a work in progress [27, 38].

To more accurately and conveniently obtain the necessary knowledge, we utilize language models as knowledge sources, prompting them to generate information. Large language models have encoded knowledge from a variety of databases into their parameters after being pre-trained on various corpora [14]. Additionally, models that undergo reinforcement learning from human feedback (RLHF) [33] can follow human instructions, serving as a natural knowledge source and providing miscellaneous knowledge based on appropriate input instructions. Based on this, we propose to use self-generated knowledge which is fetched via prompting LLMs as such:

$$p(K) = \prod_{i=1..k} p_\theta(K_i|X, K_{<i}), K_{<1} = \emptyset \qquad (4)$$

where $k$ is the length of generated knowledge tokens. When problem descriptions contain implicit intents, such as in StackOverflow [26], there is often a gap between the detailed knowledge required and the words used to describe the problem. This gap arises because deriving the required knowledge involves reasoning. To narrow this reasoning gap and obtain more precise knowledge, we decompose this extraction process when intents are implicitly given:

$$p(K) = \prod_{i=1..k} p_\theta(K_i|c, K_{<i}) \cdot p_\theta(c|X), K_{<1} = \emptyset \qquad (5)$$

where $c$ is a trial code solution from LLM based only on problem contexts. It contains necessary but potentially misused knowledge that can benefit the extraction of knowledge ($K$). $K$ can be formatted as any problem-specific structure to fit the problem instruction $X$, making it suitable for various tasks. When problem descriptions $X$ contain explicit intents, SELFEVOLVE uses Eq. 4 instead, as LLM can easily extract knowledge with high accuracy in this case.

**Revision of generated solution** Previous studies have shown that intermediate results generated from LLM may contain mistakes [29, 46, 48, 54]. Such errors can introduce noise to the prompt context, reducing the accuracy of the final output. To reduce code errors, we mimic the debugging process of programmers and introduce an iterative self-refinement mechanism to rectify the buggy code. This mechanism leverages an external Python interpreter to revise erroneous programs. Our approach incorporates code context and sample test cases into the input prompt, along with the generated code solution, to form an executable program. We then execute this program in a sandbox environment to receive error information as well as standard output. Once error information is obtained, we prompt language models to revise the buggy programs, conditioned both on program requirements and error information:

$$P(Y'|X, Y, K, e) = p_\theta(Y'|X, Y, e) \cdot p_\theta(Y|X, K) \tag{6}$$

The revised output, $Y'$, may still contain bugs. Therefore, the above process is repeated until the code can be interpreted without exceptions, or until the iteration steps reach a fixed threshold. In practical applications, the modeling of $p_\theta(Y'|X, Y, e)$ varies depending on the type of error. For simplicity, SELFEVOLVE only corrects API errors and incorrect assertion statements. We find that correcting these two types of errors contributes significantly to performance improvement in the empirical experiments.

In conclusion, we combine two LLM-driven methods - generation based on self-generated knowledge and refinement via error message - to create a more effective method called SELFEVOLVE. These two components reinforce each other almost orthogonally. On one hand, self-generated knowledge boosts self-refinement steps. By conditioning on knowledge from models' parameters, intermediate output explicitly applies the knowledge. With more accurate output, the self-refinement steps require fewer iterations to repair the code, resulting in lower difficulty. On the other hand, the self-refinement steps improve the application of generated knowledge. The input knowledge may contain irrelevant information or noise during the generation process. The self-refinement steps eliminate this noise by introducing an external interpreter, improving the overall quality of the generated code. Later empirical experiments will demonstrate how these two modules reinforce each other.

# 4 Experiments

In this work, we present a novel pipeline that supports natural and reliable code generation for a variety of programming and data science problems. To evaluate its effectiveness, we conducted experiments using three different code generation tasks: data science code generation, simple algorithm coding, and C++-to-Python code translation. These tasks were assessed using the DS-1000 [26], HumanEval [9], and TransCoder [43] benchmarks, respectively. In all experiments, we set the top-p cutoff to 0.95 and the maximum generation length to 1024. For the specific prompts used for each task, please refer to Appendix B and C.

## 4.1 Baselines

1. **DocPrompting** [56]: DocPrompting improves the LLM by retrieving problem-relevant documentation via a finetuned retriever, then condition on those documents and problem description to generate code. We use the same documentation pool for DocPrompting in DS-1000, as the problem source of DS-1000 is the same as that of CoNaLa [52]. We also use the same retrieval weights released by them, as DS-1000 are also built to test Python programming.
2. **Self-Debugging** [10] relies on a SQL application or Python interpreter to teach language models to revise SQL commands or Python code with bugs. They propose three debugging ways, including "simple", "unit test" and "explanation". Without training sets in DS-1000 and HumanEval, we implement it in a zero-shot way and use its "simple" variant for a fair comparison.
3. **SELFEVOLVE**: SELFEVOLVE is the code generation pipeline proposed in this work. In the main experiments, We use ChatGPT as the knowledge generator and the code refiner.

## 4.2 Main Results of SELFEVOLVE

**Data science code generation** For data science code generation tasks, We selected the DS-1000 [26] benchmark, which contains 1000 problems covering seven common data science libraries. DS-1000 introduces a novel "perturbation" concept, including Origin, Surface, Semantic, and Diff-Rewrite,

Table 1: Pass@1 results on the DS-1000 dataset. $^\dagger$ denotes that the results are referred from [26]. Other baselines are implemented with the same prompt and hyperparameter setting.

| Method | Perturbation | | | | Overall |
| --- | --- | --- | --- | --- | --- |
| | Origin | Surface | Semantic | Diff-Rewrite | |
| *Prior work* | | | | | |
| Codex (Completion)$^\dagger$ | 44.93 | 37.94 | 34.35 | 16.94 | 39.20 |
| Codex (Insertion)$^\dagger$ | 47.76 | 50.18 | 38.39 | 21.05 | 43.30 |
| DocPrompting | 53.95 | 50.00 | 39.57 | 25.93 | 45.50 |
| Self-Debugging | 63.38 | 59.21 | 45.65 | 28.40 | 53.00 |
| *This work* | | | | | |
| ChatGPT | 60.31 | 52.63 | 41.30 | 26.54 | 49.30 |
| SELFEVOLVE | **66.23** | **67.11** | **48.70** | **33.95** | **57.10** |
| *w/o self-refinement* | 60.09 | 59.21 | 41.30 | 29.01 | 50.60 |

representing the difficulty of problems in ascending order, making it a challenging benchmark. In this study, we prompted language models to generate problem-relevant API documentation as domain-required knowledge. For the self-refinement module, we checked the executable programs and prompted language models to fix syntax errors only. We used greedy decoding and reported the pass@1 [9] score for each method. Results are presented in Table 1. Without further refinement steps, SELFEVOLVE has already exceeded the strong ChatGPT baseline on the Surface and Diff-Rewrite perturbation types, by a margin of 6.58 and 2.47, respectively. Moreover, with an additional self-debug module, SELFEVOLVE substantially improves over other baselines, with a 7.8 (*relatively 15.8%*) pass@1 gain compared to ChatGPT, on average. SELFEVOLVE also surpasses the prompt-based method, Self-Debugging, by a convincing 4.1 performance margin. We also noticed that integrating the self-generated knowledge with the self-refinement module gains much higher improvement. Specifically, SELFEVOLVE improves the baseline in terms of all perturbation types, demonstrating that our method can impressively enhance the robustness of large language models.

**General code generation**   For general code generation tasks, we evaluated SELFEVOLVE on HumanEval [9]. This dataset contains 164 hand-written Python programming problems with an average of 7.7 test cases each. We implemented Self-Debugging methods on ChatGPT and reported its score. We did not implement DocPrompting since no library documentation is required in HumanEval. We also introduced the GPT-4 results from [5] for comparison. We reported a pass@1 score for greedy decoding and pass@10 for 10-sample decoding. For the 10-sample generation, we conducted a grid search to set the temperature to $t = 1$. Unlike DS-1000, we induced LLM to explicitly output problem-related algorithms as external knowledge and taught LLM to fix assertion errors and syntax errors. The results in Table 2 demonstrate that the strong ChatGPT baseline significantly benefits from our SELFEVOLVE method, with an 11.59 pass@1 gain and a 6.71 pass@10 gain. This leaves a small 3.95 pass@1 gap from GPT-4. Notably, with self-generated knowledge, the self-refinement module again harvested a larger improvement (+3.66 pass@10) than only applying a refinement module like Self-Debugging (+1.22 pass@10). This empirically verifies that self-generated knowledge helps to reduce most errors and produce more precise results.

**Python code translation**   As suggested by Roziere et al. [43], we experimented with our methods on the TransCoder [43] dataset. We used its test set, which requires translating C++ code to Python, and filtered out problems without testing scripts, resulting in 410 valid problems. In addition to pass@1, we followed Roziere et al. [43] by using another evaluation metric, computational accuracy, to test each model. Computational accuracy computes the accuracy score for each problem in a competition rule, where each sample code is scored as the percentile of its passed test cases, while the pass@1 metric is computed as whether the sample code has passed all test cases. We prompted LLM to generate the algorithm detail of the C++ code, which serves as the context for Python code generation. Results in Table 3 indicate that our proposed method, SELFEVOLVE, achieves the best performance among other prompt-based methods, even outperforming Self-Debugging which peeps one ground truth test case. Built upon a strong ChatGPT baseline, SELFEVOLVE further improves

Table 2: Pass@1 and pass@10 scores comparisons with different methods on HumanEval. We use the same prompt to implement each method. † denotes that scores are cited from [5].

| Model | Pass@1 | Pass@10 |
|---|---|---|
| *Prior Work* | | |
| GPT-4† | 82.00 | - |
| text-davinci-003† | 65.00 | - |
| ChatGPT | 66.46 | 86.58 |
| CodeT [8] | 65.20 | 86.80 |
| Self-Debugging | 73.78 | 87.80 |
| *Ours* | | |
| SELFEVOLVE | **78.05** | **93.29** |
| *w/o self-refinement* | 70.73 | 89.63 |

Table 3: Performance comparison on TransCoder dataset where we follow [10, 43] to translate C++ code to Python code. All methods in this work are implemented with greedy decode. "Acc." refers to computational accuracy.

| Method | TransCoder | |
|---|---|---|
| | Acc. | Pass@1 |
| *Piror Work* | | |
| PaLM [11] | 51.8 | - |
| PaLM-Coder [11] | 55.1 | - |
| Codex [9] | 80.4 | - |
| Self-Debugging [10] | 89.3 | - |
| *Ours* | | |
| ChatGPT | 92.7 | 90.0 |
| SELFEVOLVE | **94.8** | **92.4** |
| *w/o self-refinement* | 93.4 | 90.5 |

by 2.1 computational accuracy and 2.4 pass@1, respectively. Without the self-refinement module, SELFEVOLVE still improves over ChatGPT, with a 0.7 accuracy gain and 0.5 pass@1 gain.

## 4.3 Discussion

In this section, we conduct various analysis experiments to validate the efficacy of our proposed SELFEVOLVE . We first present the impact of the number of iteration step on the final performance of SELFEVOLVE. After that, we demonstrate how our generated knowledge is superior to retrieved knowledge, through a human evaluation experiment. Finally, we extend our method to an even more intelligent language model (GPT-4) to empirically show the scalability of SELFEVOLVE.

**How do iteration steps affect performance?** In §3.2, we declared that the refinement module is iteratively run to fix bugs. In this experiment, we determine under what conditions we should stop the refinement module. We ran the self-refinement module for different iteration turns on three datasets using greedy decoding and tested the pass@1 score of ChatGPT using the same prompt for each iteration stage. Figure 2a presents the detailed results. We observed that the major improvement came from the first refinement step for the HumanEval and TransCoder datasets. On DS-1000, however, the performance improved almost uniformly as we increased the number of refinement steps until the third iteration. This discrepancy across datasets resulted from the much more difficult nature of the problems in DS-1000 compared to the other two datasets. Therefore, the self-refinement



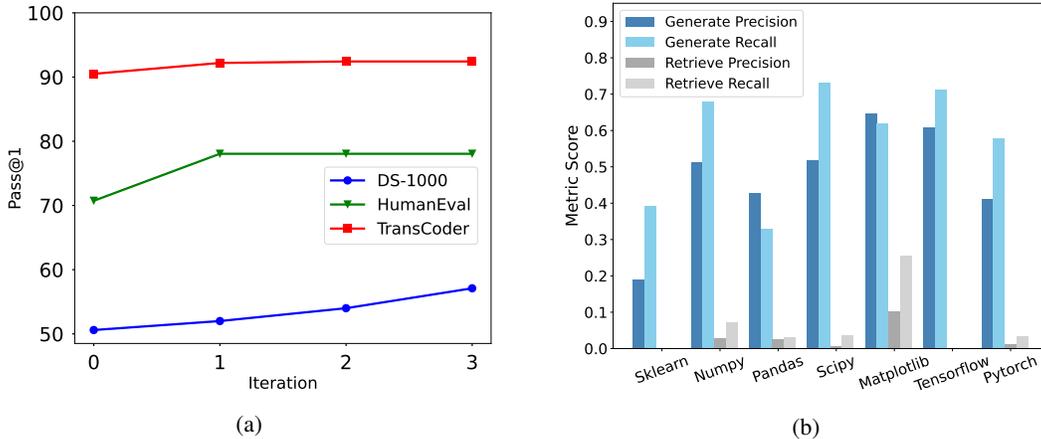(a)                                                (b)

Figure 2: (a) Performance-iteration curves of SELFEVOLVE on DS-1000, HumanEval and TransCoder datasets. (b) Precision and recall comparisons between generated knowledge and retrieved one.

Table 4: Comparison between SELFEVOLVE using ChatGPT and GPT-4 baselines. We bind SELFE-VOLVE with ChatGPT and GPT-4 to test its generalization.

| Method | DS-1000 | | | | HumanEval |
| | Scipy | Pytorch | Sklearn | Matplotlib | |
|---|---|---|---|---|---|
| SELFEVOLVE (ChatGPT) | 52.83 | 64.71 | 73.04 | 78.06 | 78.05 |
| GPT-4 | 52.83 | 44.12 | 60.00 | 69.03 | 82.00 |
| SELFEVOLVE (GPT-4) | **58.49** | **70.59** | **70.43** | **84.52** | **89.02** |
| *w/o self-refinement* | 57.55 | 63.24 | 66.09 | 69.03 | 84.76 |

module brought consistent improvement under different refinement stages. This finding suggests that SELFEVOLVE requires more refinement steps when processing difficult problems, but only one debugging turn is sufficient to bring major improvement for less complicated problems.

**Human evaluation of self-generated knowledge**    To better understand the superiority of generated knowledge in realistic scenarios, we conducted a human evaluation study to demonstrate that the generated knowledge is more relevant to the problem topic than a retrieval-based one. We randomly selected 200 problems from seven libraries of DS-1000 and asked two data science experts to count the number of correctly provided API documents to determine whether the API knowledge matched the solution. We then used two common metrics, precision and recall [6], to assess the accuracy of the knowledge in accordance with the oracle answer. Precision is defined as the percentage of correctly provided documents in the provided document set, while recall measures the percentage of correct document items in the oracle document set. More details on the experiment are presented in Appendix A. The comparison bar chart is shown in Figure 2b. We observed that in all libraries, the generated knowledge was much more accurate than the retrieved one in both metrics. Notably, the retrieved knowledge showed little match with oracle solutions in most libraries because the retrieval queries in DS-1000 are too complicated and contain implicit API demands. In the Matplotlib library, where the queries are simple and the demands are explicitly stated, the retrieved knowledge matched the problem requirements slightly but still lagged far behind the generated one. One key reason for the superiority of generated knowledge is that LLMs can bridge the reasoning gap between problem descriptions and knowledge terminology better than a retriever model. This is also the reason behind Eq. 5. In other words, generated knowledge is able to provide a more comprehensive and accurate understanding of the problem topic, which is crucial in realistic scenarios.

**Scaling to more powerful models**    To evaluate the scalability of SELFEVOLVE in more advanced language models, we integrated SELFEVOLVE with GPT-4 without requiring excessive prompt engineering. GPT-4 has demonstrated significantly greater intelligence and reasoning abilities compared to ChatGPT [5]. However, due to limited GPT-4 API-Key access, we only conducted experiments on the Scipy, Pytorch, Sklearn, and Matplotlib libraries of DS-1000, which includes a total of 444 problems and HumanEval. We used the same prompt as ChatGPT and used greedy decoding to report the pass@1 score. The results for both datasets are shown in Table 4. This experiment demonstrates that our approach can benefit from more advanced backbone models instead of degrading them. In contrast to the ChatGPT-based version, SELFEVOLVE on GPT-4 achieved higher pass@1 scores on Scipy (+4.72), Pytorch (+19.12), and Sklearn (+6.09). A similar performance gain was also observed in HumanEval, where SELFEVOLVE improved the already higher pass@1 score by 2.76. Furthermore, after adding the lightweight self-refinement plugin, GPT-4 demonstrated further improvements on all datasets. This highlights the effectiveness of leveraging advanced backbone models and the potential of our approach in producing superior results.

## 4.4   Case Study

This section demonstrates the effectiveness of SELFEVOLVE with two representative examples shown in Figure 3. In the first example, LLM generates specific documentation `tf.one_hot` for the problem. Compared to the vanilla output without documentation, language models output extra codes that cannot generalize to other test cases. In contrast, conditioning on the concrete API documentation, language models output more deterministic code without transforming the original labels to a tensor type. The provided documentation sharpens the probability curve and contributes to a more accurate
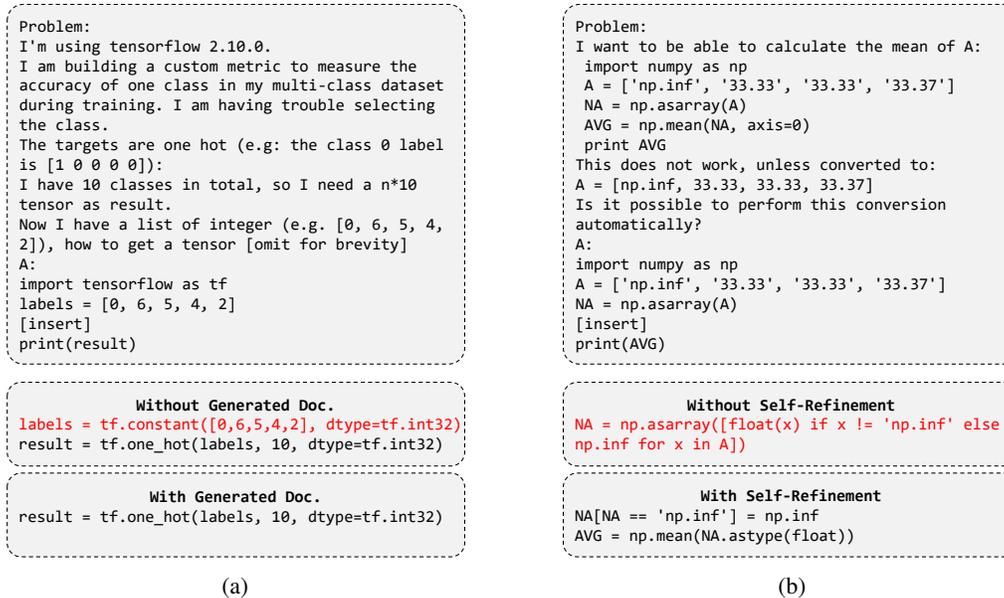
```
Problem:
I'm using tensorflow 2.10.0.
I am building a custom metric to measure the
accuracy of one class in my multi-class dataset
during training. I am having trouble selecting
the class.
The targets are one hot (e.g: the class 0 label
is [1 0 0 0]):
I have 10 classes in total, so I need a n*10
tensor as result.
Now I have a list of integer (e.g. [0, 6, 5, 4,
2]), how to get a tensor [omit for brevity]
A:
import tensorflow as tf
labels = [0, 6, 5, 4, 2]
[insert]
print(result)
```

```
Problem:
I want to be able to calculate the mean of A:
 import numpy as np
 A = ['np.inf', '33.33', '33.33', '33.37']
 NA = np.asarray(A)
 AVG = np.mean(NA, axis=0)
 print AVG
This does not work, unless converted to:
A = [np.inf, 33.33, 33.33, 33.37]
Is it possible to perform this conversion
automatically?
A:
import numpy as np
A = ['np.inf', '33.33', '33.33', '33.37']
NA = np.asarray(A)
[insert]
print(AVG)
```

```
                Without Generated Doc.
labels = tf.constant([0,6,5,4,2], dtype=tf.int32)
result = tf.one_hot(labels, 10, dtype=tf.int32)
```

```
                Without Self-Refinement
NA = np.asarray([float(x) if x != 'np.inf' else
np.inf for x in A])
```

```
                With Generated Doc.
result = tf.one_hot(labels, 10, dtype=tf.int32)
```

```
                With Self-Refinement
NA[NA == 'np.inf'] = np.inf
AVG = np.mean(NA.astype(float))
```

(a)                                    (b)

Figure 3: Two examples to show the efficacy of our proposed SELFEVOLVE methods, where red codes are wrong codes. (a) Comparison between with and without generated documentation. (b) Comparison between with and without self-refinement module.

and general answer. In the second example, language models read the `np.asarray` documentation, but forgets to compute the `AVG` value. By leveraging the traceback information without the specific test cases, language models can revise their code and use a more general method to solve the problem. Both examples illustrate how the two methods in SELFEVOLVE enhance each other to improve performance. SELFEVOLVE helps make the output of language models more general and accurate.

## 5 Limitation & Future Work

Although SELFEVOLVE has shown promising results in generating knowledge and improving the performance of large language models, there are still some limitations that need to be addressed. One of the main challenges of SELFEVOLVE is that it may not always be automatic when used in different tasks due to the hand-written prompting words. This means that its effectiveness may be limited when applied to other use cases. Another limitation of SELFEVOLVE is that the generated knowledge may not always be suitable for every task and may require fine-grained selection to be effective. However, these issues can be mitigated by developing suitable prompting skills. For example, a more comprehensive set of prompting words could be developed to make it easier to adapt SELFEVOLVE to new tasks. Additionally, a more sophisticated algorithm could be developed to automatically select the appropriate knowledge for a given task. We believe that addressing these issues will make SELFEVOLVE a more versatile and useful framework in different contexts.

## 6 Conclusion

We propose SELFEVOLVE, a simple yet effective method for solving code generation problems using large language models (LLMs) as a fully LLM-driven framework. It acts as both a knowledge provider and a self-reflective programmer to generate high-quality code in two steps, both of which are run with a single LLM. This makes it more flexible and extendable to other datasets, such as Spider [53] or APPS [19], without requiring an extra retriever or previously set up database. Substantial experiments on diverse code generation tasks have verified that SELFEVOLVE can bring great performance gains under various tasks and datasets, and outperform two strong prompting-based methods by a good margin. Furthermore, various analysis experiments indicate that SELFEVOLVE is more suitable for providing problem-related knowledge compared to traditional dense retrievers, and can be easily scaled to more intelligent language models to bring further improvements.

# References

[1] Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. A review on language models as knowledge bases. *arXiv preprint arXiv:2204.06031*, 2022.

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

[3] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR, 2022.

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[5] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023.

[6] Michael Buckland and Fredric Gey. The relationship between recall and precision. *Journal of the American society for information science*, 45(1):12–19, 1994.

[7] Shubham Chandel, Colin B Clement, Guillermo Serrato, and Neel Sundaresan. Training and evaluating a jupyter notebook data science assistant. *arXiv preprint arXiv:2201.12901*, 2022.

[8] Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. Codet: Code generation with generated tests. *CoRR*, abs/2207.10397, 2022. URL https://doi.org/10.48550/arXiv.2207.10397.

[9] 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.

[10] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.

[11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

[12] Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen tau 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://doi.org/10.48550/arXiv.2204.05999.

[13] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL https://aclanthology.org/2021.emnlp-main.552.

[14] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.

[15] Kavi Gupta, Peter Ebert Christensen, Xinyun Chen, and Dawn Song. Synthesize, execute and debug: Learning to repair for neural program synthesis. *Advances in Neural Information Processing Systems*, 33: 17685–17695, 2020.

[16] Rahul Gupta, Soham Pal, Aditya Kanade, and Shirish Shevade. Deepfix: Fixing common c language errors by deep learning. In *Proceedings of the aaai conference on artificial intelligence*, volume 31, 2017.

[17] Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. Language models can teach themselves to program better. *CoRR*, abs/2207.14502, 2022. URL https://doi.org/10.48550/arXiv.2207.14502.

[18] Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. Language models can teach themselves to program better, 2023.

[19] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.

[20] Junjie Huang, Chenglong Wang, Jipeng Zhang, Cong Yan, Haotian Cui, Jeevana Priya Inala, Colin Clement, Nan Duan, and Jianfeng Gao. Execution-based evaluation for data science code generation models. *arXiv preprint arXiv:2211.09374*, 2022.

[21] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. Mapping language to code in programmatic context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1643–1652, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1192. URL https://aclanthology.org/D18-1192.

[22] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv*, 2208, 2022.

[23] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 2022.

[24] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.

[25] Matthew Jin, Syed Shahriar, Michele Tufano, Xin Shi, Shuai Lu, Neel Sundaresan, and Alexey Svyatkovskiy. Inferfix: End-to-end program repair with llms. *arXiv preprint arXiv:2303.07263*, 2023.

[26] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science code generation. *arXiv preprint arXiv:2211.11501*, 2022.

[27] Mukai Li, Shansan Gong, Jiangtao Feng, Yiheng Xu, Jun Zhang, Zhiyong Wu, and Lingpeng Kong. In-context learning with many demonstration examples. *arXiv preprint arXiv:2302.04931*, 2023.

[28] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.

[29] Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2301.13379*, 2023.

[30] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023.

[31] Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. Augmented language models: a survey. *arXiv preprint arXiv:2302.07842*, 2023.

[32] 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_.

[33] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[34] Chris Parnin and Alessandro Orso. Are automated debugging techniques actually helping programmers? In *Proceedings of the 2011 international symposium on software testing and analysis*, pages 199–209, 2011.

[35] Md Rizwan Parvez, Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. Retrieval augmented code generation and summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2719–2734, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.232. URL https://aclanthology.org/2021.findings-emnlp.232.

[36] Panupong Pasupat, Yuan Zhang, and Kelvin Guu. Controllable semantic parsing via retrieval augmentation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7683–7698, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.607. URL https://aclanthology.org/2021.emnlp-main.607.

[37] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*, 2023.

[38] Bo PENG. RWKV-LM. https://github.com/BlinkDL/RWKV-LM, 8 2021.

[39] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1250. URL https://aclanthology.org/D19-1250.

[40] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL https://arxiv.org/abs/1908.10084.

[41] Stephen E Robertson and Steve Walker. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*, pages 232–241. Springer, 1994.

[42] Tobias Roehm, Rebecca Tiarks, Rainer Koschke, and Walid Maalej. How do professional developers comprehend software? In *2012 34th International Conference on Software Engineering (ICSE)*, pages 255–265. IEEE, 2012.

[43] Baptiste Roziere, Marie-Anne Lachaux, Lowik Chanussot, and Guillaume Lample. Unsupervised translation of programming languages. *Advances in Neural Information Processing Systems*, 33, 2020.

[44] Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I. Wang. Natural language to code translation with execution. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3533–3546, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.emnlp-main.231.

[45] Ke Wang, Rishabh Singh, and Zhendong Su. Dynamic neural program embeddings for program repair. In *International Conference on Learning Representations*, 2018.

[46] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.

[47] Yue Wang, Weishi Wang, Shafiq Joty, and Steven C.H. 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*, pages 8696–8708, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.685. URL https://aclanthology.org/2021.emnlp-main.685.

[48] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.

[49] Frank F. Xu, Zhengbao Jiang, Pengcheng Yin, Bogdan Vasilescu, and Graham Neubig. Incorporating external knowledge through pre-training for natural language to code generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6045–6052, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.538. URL https://aclanthology.org/2020.acl-main.538.

[50] Michihiro Yasunaga and Percy Liang. Graph-based, self-supervised program repair from diagnostic feedback. In *International Conference on Machine Learning*, pages 10799–10808. PMLR, 2020.

[51] Michihiro Yasunaga and Percy Liang. Break-it-fix-it: Unsupervised learning for program repair. In *International Conference on Machine Learning*, pages 11941–11952. PMLR, 2021.

[52] Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. Learning to mine aligned code and natural language pairs from stack overflow. In *International Conference on Mining Software Repositories*, MSR, pages 476–486. ACM, 2018. doi: https://doi.org/10.1145/3196398.3196408.

[53] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3911–3921, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1425. URL https://aclanthology.org/D18-1425.

[54] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.

[55] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR, 2021.

[56] Shuyan Zhou, Uri Alon, Frank F. Xu, Zhengbao Jiang, and Graham Neubig. Docprompting: Generating code by retrieving the docs. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=ZTCxT2t2Ru.

# A  Comparison between Generated Knowledge and Retrieved Knowledge

We here present the experiment details of this human evaluation experiment. We randomly select one-fifth of the problems in each library to form the 200 problem set, whose formation is shown in Table 5. API is defined as a function call, a method call or an attribute getter and setter in this case. For retrieved knowledge, we follow DocPrompting [56] to take problem descriptions as queries and their provided document pool as the target sets to perform the retrieval. We use their pretrained CodeT5 retriever and retrieve $k = 5$ knowledge items from the target pool. For irrelevant documents, we filter them for both retrieved documents and generated ones, so the final number of documents for each problem may be smaller than $k$. After retrieval, we follow Zhao et al. [55] to put items with higher scores near the generation position to ensure the best generation result.

Table 5: Problem counts for each library in the selected set.

| Library | Tensorflow | Pytorch | Numpy | Matplotlib | Pandas | Sklearn | Scipy |
|---------|-----------|---------|-------|-----------|--------|---------|-------|
| #Problems | 10 | 15 | 40 | 37 | 44 | 28 | 26 |

# B  Prompt for the First-Stage of SELFEVOLVE in Each task

We show the detailed prompt for the first step of SELFEVOLVE used in each task below. We show prompts for DS-1000, HumanEval and TransCoder in Figure 4, Figure 5, and Figure 6, respectively.

# C  Prompt for Self-Refinement in Each task

We show the detailed prompt for the self-refinement module used in each task below. We show prompts for DS-1000, HumanEval and TransCoder in Figure 7, Figure 8, and Figure 9, respectively.

```
problem:                                                          Code Snippet
{problem description}

A:
{code context}
[insert]
print(result)

Could you help to write solution codes and store the answer in the variable result?
You only need to output the codes which can fill the [insert] block.
Please just output codes without any explanation and natural language.
Wrap your code with "```".
```

```
This is a snippet of {library} code:                             api list
{code_snippet}


I want you to show the API used in the code.
Here are some rules for displaying the APIs:
1. do not show basic API like print(), import, __str__, __rpre__, etc. Just show API about
{library}
2. for class methods, make sure add the class name and library name before method with only
two dots separated, like tf.random.normal.
3. do not output duplicate APIs

Here's what APIs the above code calls line by line:
```

```
Please show me the following API's API specification:              generated doc
{api list}

You do not need to show me examples of each API.
Your answer should start with "1. "
```

```
Documentation                                                     final solution
{generated doc}

problem:
{problem description}

A:
{code context}
[insert]
print(result)

Could you help to write solution codes and store the answer in the variable result?
You only need to output the codes which can fill the [insert] block.
Please just output codes without any explanation and natural language.
Wrap your code with "```".
```

Figure 4: Prompt for the first step of SELFEVOLVE in the DS-1000 dataset.

```
                                                                    problem-related algorithms
```
{Problem description}
```

For the above question, could you briefly teach me how to solve it step by step in natural
language?
Don' t write the code in this step.
```

```
                                                                    final solution
Based on the above idea, help me complete the function.
Be attention, you should only output the codes without any explanation and natural language.
Wrap your code with "```"
```

Figure 5: Prompt for the first step of SELFEVOLVE in the HumanEval dataset.

```
Can you help me to explain how a piece of code work from the perspective of algorithms and
syntaxes? here is the code:

```
{Cpp code}
```

Your explanation:                                                   explanation
```

```
Based on your explanation, write a same function in python with the same function name, the
same function argument and the same functionality.
                                                                    final solution
```

Figure 6: Prompt for the first step of SELFEVOLVE in the TransCoder dataset.

```
{problem description}                                               explanation
{code context}
{error code}
After running the above code, it raises such error:
{error message}

It seems that this line
{bug_code}
has bugs.
Can you tell me what causes this bug?
```

```
                                                                    final solution
Conclude your code so that your answer should fit in the following code context:
{error code}

You only need to output codes that can fill in the [insert] block.
You do not need to output the codes before the [insert] block.
You only need to output the codes without explanation.
Wrap your code with "```".
```

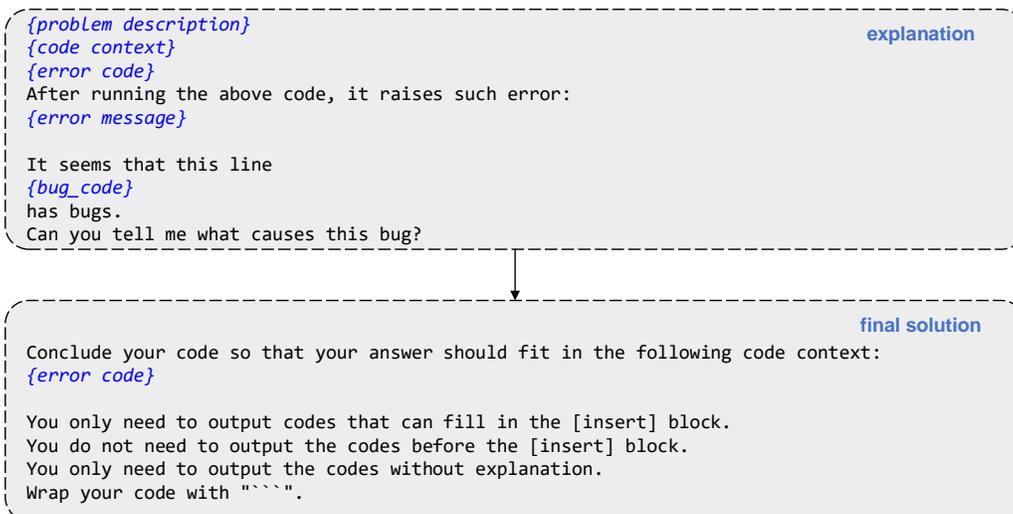Figure 7: Prompt for self-refinement in the DS-1000 dataset.

16

```
                                                                    revise syntax
{error code}
```

When I run {test case}, I meet {syntax error}.
Help me refine the code.
You should only output the codes without any explanation and natural language.
Wrap your code with "```"

```
                                                                    revise error
{error code}
```

The expected output of {test case} is {ground truth}. However, the above code output {wrong output}.
Help me refine the code.
You should only output the codes without any explanation and natural language.
Wrap your code with "```"

Figure 8: Prompt for self-refinement in the HumanEval dataset.

```
                                                                    revise syntax
{error code}

After running the above code, it raises such error:
{error_msg}

Can you tell me how to fix this bug?
You only need to output the codes.
You do not need to output explanations.
Wrap your code with "```" and "```".
```

```
                                                                    revise error
I have a c++ function like this:
```
{cpp_code}
```

I write a python function to implement the same functionality with the above c++ code:

```
{py_code}
```

But the python function has some bugs and I cannot find them.
Can you tell me which part of my python code is different from C++ code?
Help me correct those codes.
Wrap your code with "```".
```

Figure 9: Prompt for self-refinement in the TransCoder dataset.