

~~DF1.APPEND(DF2)~~ PD.CONCAT([DF1,DF2])

GITCHAMELEON: UNMASKING THE VERSION-SWITCHING CAPABILITIES OF CODE GENERATION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid evolution of software libraries presents a significant challenge for code generation models, which must adapt to frequent version updates while maintaining compatibility with previous versions. Existing code completion benchmarks often overlook this dynamic aspect, and the one that does consider it relies on static code prediction tasks without execution-based evaluation, offering a limited perspective on a model’s practical usability. To address this gap, we introduce **GitChameleon**, a novel, manually curated dataset comprising 116 Python code completion problems, each conditioned on specific library versions and accompanied by executable unit tests. **GitChameleon** is designed to rigorously assess the ability of modern large language models (LLMs) to generate version-specific code that is not only syntactically correct but also functionally accurate upon execution. Our comprehensive evaluations reveal that state-of-the-art LLMs struggle with this task; for instance, **GPT-4** achieves a pass@10 of only 39.9% (43.7% when provided with error feedback), highlighting the complexity of the problem and the limitations of current models. By providing an execution-based benchmark that emphasizes the dynamic nature of code libraries, **GitChameleon** serves as a critical tool for advancing the development of more adaptable and reliable code generation models. We release the dataset and evaluation framework to encourage further research in this vital area.

1 INTRODUCTION

Large Language Models (LLMs) have become highly popular in code completion, to the extent that they are now deployed as virtual coding assistants within popular code editors¹, enhancing the overall coding workflow. Code, being a dynamic and constantly evolving environment, necessitates a continuous process of adaptation to stay in sync with the rapidly shifting paradigms, frameworks, and methodologies within the software development domain. The inherent variability in coding styles, the emergence of new programming languages, and the continuous evolution of libraries and packages underscore the imperative for an active approach in updating code generation models.

In response to the needs of practical coding environments, several large language models (LLMs) have been introduced, including StarCoder (Li et al., 2023), DeepSeek-Coder (Guo et al., 2024), CodeLlama (Rozière et al., 2023), among others. Despite these advancements, existing LLMs often struggle to keep pace with the rapid changes in codebases, particularly when tasked with generating version-specific code that is both syntactically and functionally accurate. This issue is especially critical, as developers increasingly depend on AI-assisted coding tools to boost productivity and maintain code quality. A recent Stack Overflow survey revealed that 70% of the participants are using or planning to integrate AI coding tools, 33% citing increased productivity as the primary motivation to incorporate these tools into their workflows².

¹<https://github.com/features/copilot>

²<https://stackoverflow.co/labs/developer-sentiment-ai-ml/>

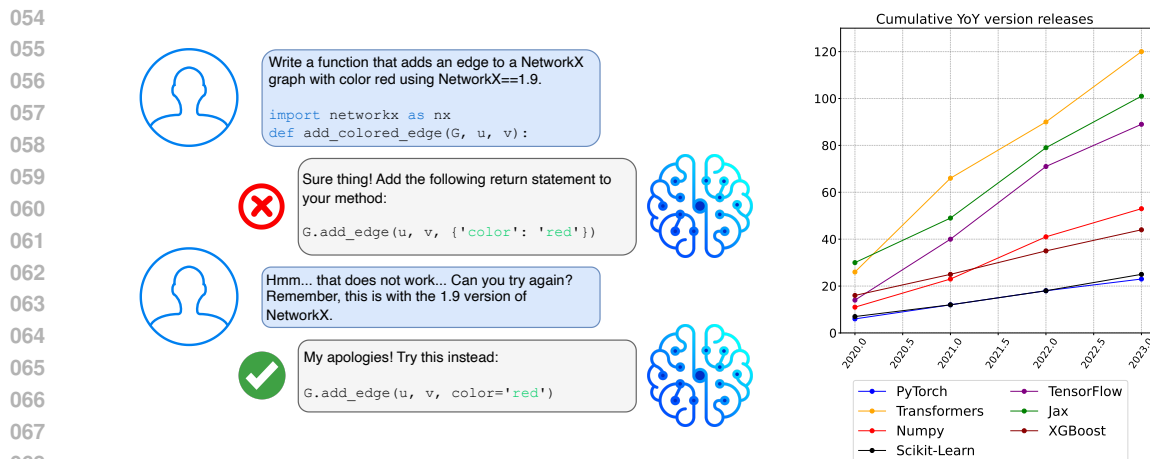


Figure 1: **Left:** Modern LLMs often struggle with generating version-accurate code, highlighting the need for benchmarks that specifically assess their ability to handle versioning. **Right:** Cumulative year-over-year version releases of popular Python-based machine learning libraries show a consistent upward trend, reflecting the rapid pace of development and version updates of code libraries and packages.

Given the rapid development and release cycles of popular libraries, as shown in Figure 1 (right), the need for code generation models to continually adapt to changing API’s is more pressing than ever. For example, prominent machine learning and deep learning libraries like `PyTorch` (Paszke et al., 2019), `NumPy` (Harris et al., 2020), and `Scikit-Learn` (Buitinck et al., 2013) undergo frequent updates, which are reflected in a consistent upward trend in user downloads and version releases. This dynamic nature of code requires models that can adapt and generate code that adheres to the latest versions and practices, a need that current models often fail to meet comprehensively. In addition, certain hardware is restricted to compatibility with specific versions of commonly used packages, which adds an additional layer of complexity beyond merely updating the knowledge base of a code LLM to the latest library versions.

In response to these challenges, our work introduces a novel benchmark designed to assess the ability of LLMs to generate version-specific code. We propose **GitChameleon**, a benchmark that evaluates state-of-the-art code models by requiring them to produce executable code based on version-specific prompts. This code is then executed to verify its correctness against expected outputs. By highlighting the limitations of current models in generating accurate version-specific code, **GitChameleon** provides a structured approach to enhance these models and ensure their practical utility in real-world coding environments.

In summary, our contributions are highlighted as follows: 1) we introduce a novel code completion benchmark **GitChameleon** consisting of 116 Python-based version conditioning problems including human written unit tests; 2) we systematically analyse the version-specific performance of state-of-the-art code generation LLMs on API change type, version release year, and specific libraries. 3) we demonstrate the effectiveness of utilizing error log feedback as a way to improve version conditioning performance of code generation LLMs.

2 GITCHAMELEON BENCHMARK

We introduce **GitChameleon**, a benchmark comprising 116 Python-based version conditioning problems focused on popular code libraries. To evaluate LLM performance on **GitChameleon**, each problem is accompanied by handwritten assertion-based unit tests, enabling a thorough execution-based assessment of the outputs generated by the code LLMs. This structured approach enables a thorough understanding and categorization of LLM failures in common scenarios involving version-specific code generation problems. In the following sections, we detail the benchmark statistics, data collection methodology, and sample verification process.

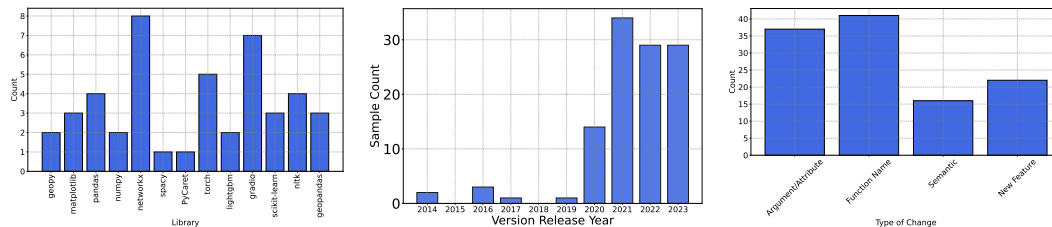
Table 1: Compared to other popular code generation benchmarks, including those evaluating version conditioning, **GitChameleon** features library- and version-specific problems with unit tests based on real version changes, closely aligning with practical settings.

Dataset	Problems	Data Source	Library Specific	Version Specific	Execution based	Real
HumanEval (Chen et al., 2021)	164	Hand-Written	✗	✗	✓	-
MBPP (Austin et al., 2021)	974	Hand-Written	✗	✗	✓	-
MTPB (Nijkamp et al., 2022)	115	Hand-Written	✗	✗	✓	-
APPS (Hendrycks et al., 2021)	10000	Competitions	✗	✗	✓	-
CodeContests (Li et al., 2022)	117	Competitions	✗	✗	✓	-
JulCe (Agashe et al., 2019)	1518049	Notebooks	✗	✗	✓	-
DSP (Chandel et al., 2022)	1119	Notebooks	✓	✗	✓	-
CoNaLa (Yin et al., 2018)	2879	StackOverflow	✓	✗	✗	-
DS-1000 (Lai et al., 2023)	1000	StackOverflow	✓	✗	✓	-
BigCodeBench (Zhuo et al., 2024)	1140	Expert-Curated	✓	✗	✓	-
Versicode (Wu et al., 2024b)	98692	GitHub, StackOverflow	✓	✓	✗	✓
CodeUpdateArena (Liu et al., 2024)	670	LLM-Generated	✓	✓	✓	✗
Wang et al. (2024)	28125	API change logs	✓	✓	✗	✓
GitChameleon (Ours)	116	Handwritten and LLM-assisted	✓	✓	✓	✓

2.1 STATISTICS

GitChameleon consists of 116 python-based version conditioned problems based on 11 libraries: PyTorch (Paszke et al., 2019), Geopandas (Jordahl et al., 2020), NLTK (Bird & Loper, 2004), NetworkX (Hagberg et al., 2008), GeoPy³, Gradio (Abid et al., 2019), Scikit-Learn (Buitinck et al., 2013), Matplotlib (Hunter, 2007), PyCaret⁴, Pandas (pandas development team, 2020; Wes McKinney, 2010) and NumPy (Harris et al., 2020). The samples were collected from version releases over a period from the year 2014 to 2023 and excludes legacy and yanked version releases.

Using the `cl100k_base` tokenizer, we analyzed the token counts of the **GitChameleon** samples. The problem statements average 20.4 tokens, and the starter code averages 47.4 tokens, leading to a combined average of 67.8 tokens per sample. Including the prompt template utilized for evaluating instruction-tuned LLMs, the total token count across all samples is 19,409 tokens.



(a) Number of unique versions per library in **GitChameleon**. (b) Number of samples by version release year. (c) Number of samples by type of change.

Figure 2: Fine-grained statistics of the **GitChameleon** benchmark. (a) The library with the most unique versions in the dataset is `networkx` with 8, whereas only 1 version of `spacy` and `PyCaret` are represented in **GitChameleon**. (b) Most versions in the dataset were released between 2021-2023, with a few versions released in earlier years. (c) The most common type of changes between versions were function name changes and argument/attribute changes, while semantic/output changes were least common.

As demonstrated in Fig. 2b, most of the samples in **GitChameleon** are from versions of libraries released in the years 2021, 2022 and 2023, with 2021 released version samples accounting for 35% of the total sample count in the benchmark. Since some of the models evaluated on **GitChameleon** have disclosed their training data cutoff dates, we have ensured that most, if not all, samples fall within the training window of these models. This approach helps to ensure that the models during their training period are likely exposed to the versions on which the samples are based. Fig. 2a shows that `NetworkX` (Hagberg et al., 2008) and `Gradio` (Abid et al., 2019) have the most versions in our

³<https://pypi.org/project/geopy/>

⁴<https://pycaret.org/>

benchmark (8 and 7, respectively). Meanwhile, PyTorch (Paszke et al., 2019) and NumPy (Harris et al., 2020) have the highest number of samples (18 and 15, respectively), together accounting for 34% of the total sample count. Further, we annotate each sample with the type of change that is classified into the following categories:

- **Argument or Attribute change:** The API call to a function, method, or class has a change in arguments (e.g. name, order, new, deprecated argument) between versions.
- **Function Name change:** The name of the API call has changed between versions (e.g. `pandas.append` to `pandas.concat`).
- **Semantics or Function Behaviour change:** The semantic / runtime behaviour of the API call changed between versions (e.g. returning a different type).
- **New feature or additional dependency based change:** A feature was introduced in a specific version, therefore, to execute the same functionality, a model using an older version should make use of an additional dependency (e.g. `torch.special` was introduced in TORCH 1.10, previously one could use NUMPY for the same).

Most samples in the GitChameleon benchmark fall under the Argument or Attribute and Function Name change category, as these are the most frequent and expected types of changes in mature and stable libraries.

Differentiating factor Several datasets examine LLM interactions with version-specific code, including Versicode (Wu et al., 2024b), CodeUpdateArena (Liu et al., 2024), and the dataset by Wang et al. (2024). While these datasets are valuable, our dataset offers a unique and complementary perspective by focusing on the real-world scenario where developers are often constrained to specific library versions due to technical debt. CodeUpdateArena investigates model adaptation to synthetic API changes, we focus our evaluation on real API changes to assess how effectively an LLM can generate code for version-specific changes of library versions that they have been trained with. In contrast, Versicode and Wang et al. (2024)’s datasets, while addressing library evolution, primarily rely on string matching for evaluation. Our approach diverges by incorporating executable tests, providing a more practical and rigorous assessment of code generation capabilities.

2.2 COLLECTION FRAMEWORK

Task Description	Expected Result
<pre># Write a function that checks if all elements in an array are true.</pre>	<code>np.all(arr)</code>
Starter Code	Assertion Test
<pre>import numpy as np def alltrue_fn(arr): return</pre>	<pre>arr = np.array([1, 1, 1, 1]) result = alltrue_fn(arr) assert result == np.all(arr)</pre>

Table 2: Example of a problem statement derived from a changelog entry from Numpy 1.25.0

The examples were manually crafted by the authors, who divided the task among themselves. We compiled a list of popular Python libraries, focusing on those with which at least one author was familiar and that had detailed changelogs documenting changes between versions. For each library, we reviewed the changelogs to identify deprecated functions, argument changes, alterations in behavior, and newly introduced functions.

For each identified change, we create a concise problem statement, write the starter code, define the expected solution, and develop an assertion test. For instance, in Table 2, we illustrate an example based on the changelog for version 1.25.0 of NumPy (Harris et al., 2020), a library for scientific computing in Python. This changelog notes that “`np.alltrue` is deprecated. Use `np.all` instead.” We used this change to craft a problem statement that tests the LLMs’ ability to recognize and adapt to version-specific updates.

Unit test and evaluation framework verification To assess the correctness of the evaluation framework of GitChameleon, we needed to verify three key aspects:

- **Compilation:** Ensure that the starter code compiles successfully.
- **Assertion unit tests:** Confirm that the assertion tests function correctly.
- **Dependencies:** Verify that all necessary external dependencies are installed, excluding the ones being tested.

We used `venv` to create and manage virtual environments for testing. This process involved installing the appropriate library version and any additional dependencies. We then combined the starter code, expected result, and the assertion test into a single script, which was executed to verify all three criteria. We provide pseudocode for our verification process in appendix A.2.

3 EMPIRICAL STUDY

We evaluate state-of-the-art large language models (LLMs) using the **GitChameleon** benchmark to assess their ability to generate version-specific, executable code. This study highlights how well current models adapt to dynamic library versions and produce functionally correct code that passes the provided unit tests.

3.1 EXPERIMENTAL SETUP

For each open-source LLM, we downloaded the corresponding Hugging Face (HF) weights and served the models using Text Generation Inference (TGI). We used a single NVIDIA 95GB H100 GPU for models with fewer than 70 billion parameters, two GPUs for models more than 70 billion parameters.

We configured the generation parameters with a `top_p` value of 0.95, `top_k` of 50, and a temperature of 0.3 for **Pass@1** and 0.8 for **Pass@10**, in addition to finding the optimal temperature for each model. The maximum number of new tokens generated was set to 256. Additionally, we enabled flash attention (Dao et al., 2022) for all models to enhance inference efficiency. A list of all the models and their cutoff dates if available is provided in appendix A.3.

3.2 EVALUATION METRICS

To comprehensively evaluate the performance of code generation models using the **GitChameleon** dataset, we employ a range of execution-based metrics. These metrics assess not only the correctness of the generated code but also its efficiency and adaptability to different versions.

Pass@k measures the proportion of problems for which at least one of the top k generated solutions passes all assertion tests. This metric provides insight into the model’s ability to generate functionally correct code. For each problem, we generate n code samples, and compute the pass at k metric by the corrected formula:

```

1 def corrected_pass_at_k(n, c, k=10):
2   if n - c < k: return 1.0
3   return 1.0 - np.prod(1.0 - k / np.arange(n - c + 1, n + 1))

```

For instruct models, we run the model’s parsed output as standalone code, and for base models, the concatenation of the starting code and model’s parsed output (completion).

Greedy refers to the standard greedy decoding method, where the most probable token from the next-token distribution is deterministically chosen. This is analogous to setting the temperature to 0.

Error Feedback adds the error log to the prompt (after executing the generated code from the model with the initial prompt). Then, the `pass@k` metric is recalculated based on the model’s generated code using the prompt with error feedback. See appendix A.1 for an example.

3.3 MAIN RESULTS

We report the performance of both base and instruct-tuned models on the **GitChameleon** benchmark in Tables 3 and 4, respectively. Our analysis reveals a strong positive correlation between

model size and performance in version-specific code generation tasks. For base models, Spearman’s rank correlation coefficients are 0.82 for Pass@1 and 0.69 for Pass@10 (both p-values <0.01), indicating that larger models generally perform better. Specifically, DeepSeek-Coder 33B achieved the highest Pass@1 score of 35.7%, highlighting its proficiency in generating correct solutions on the first attempt, while CodeLlama 34B outperformed others at Pass@10 with a score of 42.8%, demonstrating its ability to produce correct solutions given multiple attempts.

Model	Size	Pass@1 T=0.3	Pass@10 T=0.8
CodeLlama	7B	20.4±1.6	36.1±5.5
	13B	25.8±1.0	36.4±2.0
	34B	30.6±1.4	42.8±1.4
StarCoder2	3B	11.9±1.9	27.1±1.9
	7B	15.5±1.1	23.1±2.6
	15B	13.7±1.7	27.0±3.4
Llama-3	8B	22.3±1.0	32.0±2.1
	70B	27.2±3.0	41.3±2.5
Qwen2	7B	27.4±1.2	37.7±1.8
	72B	33.2±2.1	39.7±5.5
StarCoderbase	1B	13.3±1.0	20.3±1.2
	3B	15.5±1.2	26.5±1.5
	7B	20.0±0.9	31.3±4.1
	15B	16.9±1.8	30.8±2.6
StarCoder	15B	16.0±1.2	35.9±1.9
Deepseek-coder	1.3B	22.0±2.5	28.0±1.9
	6.7B	31.0±1.8	36.1±0.7
	33B	35.7±3.0	37.9±4.9

Table 3: **Base Model Performance Metrics.** Deepseek-coder-33B is the strongest model for Pass@1 (temperature 0.3), while CodeLlama-34B is the strongest model when we compute Pass@10 with an increased number of generations (20) sampled at temperature 0.8. We observe that there is a strong positive correlation between model size and performance, with Spearman’s rank correlation coefficients of 0.82 for Pass@1 and 0.69 for Pass@10.

Similarly, for instruct-tuned models, we observe Spearman’s rank correlation coefficients of 0.52 for Pass@1 and 0.70 for Pass@10 (both with p-values under 1%), confirming the positive correlation between model size and performance. Phi-3.5-MoE (16×3.8B) achieved the highest baseline Pass@1 score of 30.9% (33.6% greedy) and Pass@10 (40.5%). GPT-4o outperformed others at Pass@10 with error feedback with a score of 43.7%. Incorporating error feedback led to average improvements of 4.47% in Pass@1 and 3.51% in Pass@10 across instruct-tuned models. Additionally, Pass@10 with n=20 samples showed an average performance improvements of 10.6% for base models and 14.8% for instruct-tuned models over Pass@1 with n=5. These findings highlight that scaling up model size, utilizing error feedback, and allowing multiple solution attempts are effective strategies for enhancing the ability of LLMs in handling version-specific code generation tasks.

3.4 IN-DEPTH ANALYSIS OF FINDINGS

In this section, we delve deeper into the results obtained from our experiments, analyzing model performance across various dimensions, including model size, year of library release, the type of API changes encountered and sample difficulty.

Analysis of Performance by Release Date At the top of Figure 4, we present the year-over-year performance of a subset of the instruction-finetuned models. The average performance of all models dropped significantly from 87.7% in 2019 (not shown) to 28.2% in 2023, with intermediate values of 79.1%, 45.2%, and 21.3%. This decline is likely due to the fact that the training datasets contain more data from earlier years, underscoring the need for code LLMs to better adapt to the evolving

Model	Size (Context) / Version	Pass@1 ($T = 0.3$)			Pass@10 ($T = 0.8$)	
		Baseline	+ Error Feedback	Greedy [†]	Baseline	+ Error Feedback
StarCoder2-v0.1	15B	22.4±0.9	27.9±1.2	21.6	37.4±1.0	38.1±1.3
CodeLlama	7B	16.5±1.0	19.6±0.9	19.0	27.2±2.2	32.2±1.7
	13B	20.3±0.6	25.6±1.3	22.4	35.7±0.8	41.2±1.0
Llama-3.1	8B	15.7±0.5	20.0±0.3	16.4	28.8±1.6	35.1±1.3
Llama-3.2	1B	9.3±0.4	12.0±0.7	9.5	16.2±0.7	20.1±0.7
	3B	10.4±0.6	14.0±0.4	12.1	20.2±0.6	25.4±0.7
CodeQwen1.5-Chat	7B	20.9±0.4	25.4±0.5	21.6	40.2±0.8	42.4±0.8
Qwen2	7B	17.8±0.4	24.8±1.5	18.1	38.4±1.0	40.7±1.0
	72B	26.0±0.6	29.0±0.4	26.7	38.2±0.7	40.8±0.3
Qwen2.5-Coder	1.5B	19.7±0.9	22.9±1.2	19.8	34.1±0.4	37.6±0.4
	7B	21.2±0.4	24.0±0.7	22.4	35.4±1.2	41.3±0.9
Codestral-v0.1	22B	25.1±0.6	31.6±0.4	25.0	37.4±0.3	41.5±0.3
Yi-Chat	6B	17.4±0.4	23.2±1.0	15.5	33.6±0.9	36.6±0.8
	9B	19.9±0.6	24.8±0.7	20.7	30.6±0.5	39.1±0.3
	34B	20.8±0.5	26.3±1.0	21.6	35.4±0.5	38.4±0.8
codegemma	7B	17.8±0.7	22.6±1.0	16.4	33.9±0.6	38.0±0.5
stable-code	3B	14.6±0.7	16.3±0.9	14.7	23.9±1.4	25.9±0.9
granite-code	3B (128k)	23.6±1.1	27.0±1.4	22.4	33.7±0.3	34.8±0.5
	8B (4k)	24.8±0.5	28.4±1.0	24.1	39.3±1.2	41.2±0.6
	8B (128k)	23.4±0.6	27.7±1.0	25.9	35.5±1.7	38.8±1.1
	20B (8k)	28.7±0.5	30.0±0.8	29.3	37.0±0.9	37.3±0.5
	34B (8k)	29.6±0.9	31.4±1.0	30.2	37.3±1.0	40.9±0.8
Phi-3.5-mini	3.8B	24.2±0.8	29.9±1.2	26.7	35.2±1.5	37.4±0.9
Phi-3.5-MoE	16x3.8B	30.9±0.8	34.9±0.7	33.6	40.5±0.5	43.2±0.1
Nxcode-CQ-orpo	7B	20.8±0.7	25.0±1.1	21.6	38.9±1.0	42.4±0.7
GPT	3.5	19.6±1.0	27.2±0.8	19.8	33.3±1.0	37.8±1.7
	4o (2024-05-13)	23.6±2.7	34.1±1.2	25.0	39.9±2.4	43.7±2.4

Table 4: **Instruct Model performance metrics.** (Top) OSS models, (bottom) closed-sourced models. We observe a 4.47% and 3.51% improvement with error feedback in Pass@1 and Pass@10, respectively. Additionally, there is a strong positive correlation between model size and performance, with Spearman’s rank correlation coefficients of 0.52 for Pass@1 and 0.70 for Pass@10.

nature of code libraries and their versions. Interestingly, many models appear to improve with error feedback disproportionately across versions released in the years 2021-2023. For example, Qwen2-72B and Llama-3.2 3B improve more with feedback in 2022 compared to 2021 or 2023, while GPT-4o improves more with feedback in 2023. This raises a question about the extent to which models’ training data influences the effectiveness of error feedback.

Analysis of Performance by Type of API Change At the bottom of Figure 4, our analysis of model performance across different types of API changes in **GitChameleon** revealed significant variations. The models struggled the most with **Semantics or Function Behaviour** changes, achieving an average Pass@1 score of only 7.34%. **Argument and Attribute** changes were the second most challenging, with an average Pass@1 score of 18.5%. In contrast, the models performed better on **Function Name** changes and **New Feature or additional dependency based** changes, with average Pass@1 scores of 50.5% and 48.6%, respectively.

In general, larger models are more robust to the name changes, argument/attribute changes, and new feature. However, all models perform very poorly on semantic changes, regardless of the availability of error feedback. This indicates a weakness of SotA code generation models, and an area for further investigation. Furthermore, error feedback appears to have a more significant impact in argument/attribute changes compared to the other types of changes. This indicates that the models

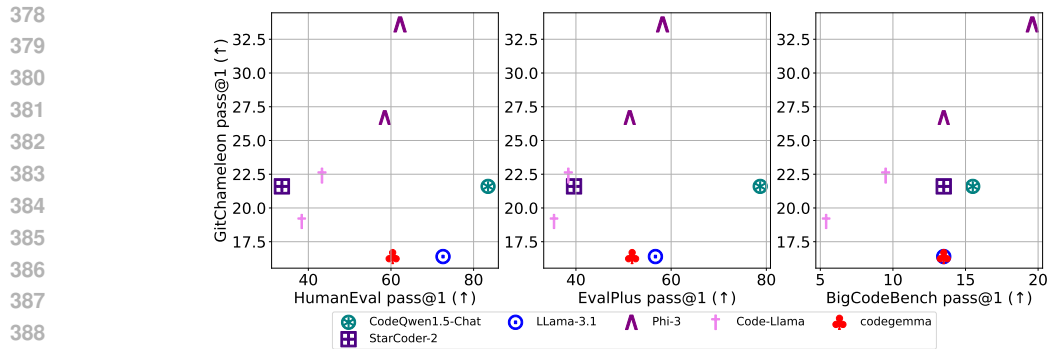


Figure 3: Correlation of **GitChameleon** with the representative code benchmarks (HumanEval, EvalPlus, and BigCodeBench-Hard Complete split). For each benchmark, the spearman correlation coefficient was -0.08 , 0.07 , and 0.35 , respectively. While HumanEval and EvalPlus showed very weak correlations, BigCodeBench-Hard showed a positive correlation ($+0.35$) with **GitChameleon**.

may be using the error feedback to directly address failures in version-conditioned code generation, rather than non-specific errors such as syntax errors.

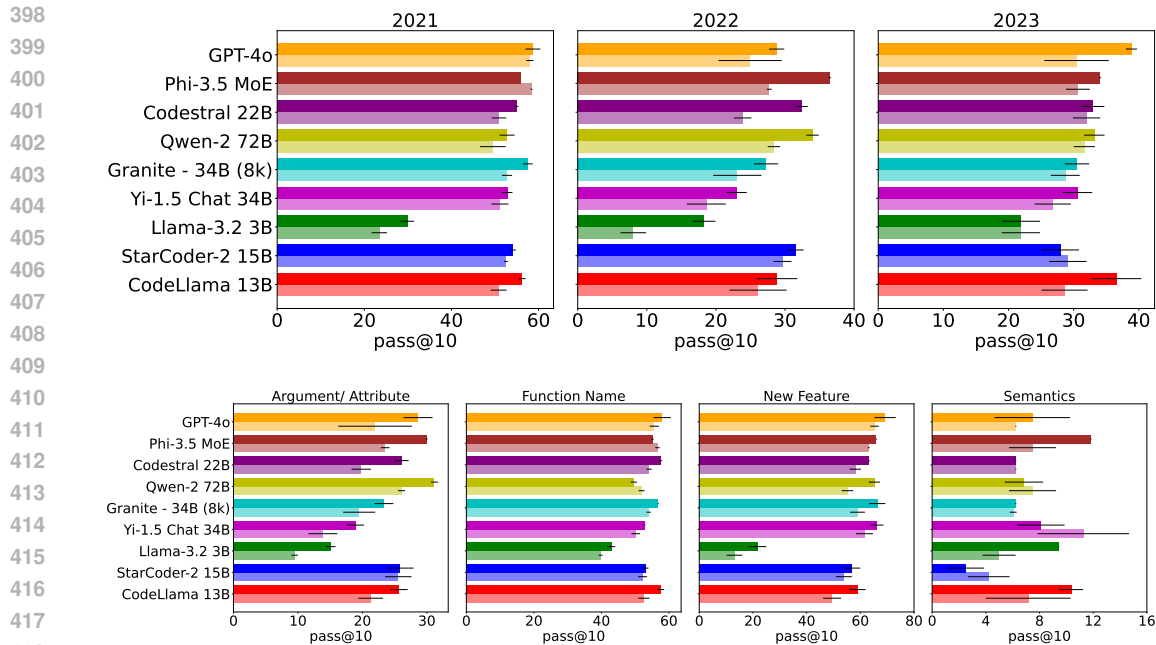


Figure 4: **Instruct-tuned model performance breakdown by version release year (top) and type of change (bottom)**: We analyze model performance in terms of $\text{pass} @ 10$ for baseline and with error feedback generation across two dimensions: version release year and type of changes. Darker shaded bars represent values obtained via error feedback generation. Standard deviation is drawn as a black line, obtained from 5 random seeds. (Top) Many models perform poorly on 2022, and generally perform worse on more recent versions. (Bottom) All models perform very poorly at semantic changes, indicating an potential area for massive improvement. Most models perform well on function name changes and new feature (with the exception of Llama-3.2-3B, which is also the smallest model in this comparison).

Sample difficulty analysis Figure 3.4 shows the distribution of sample difficulty. Notably, individual models (right panel) often display bimodal distributions, meaning they tend to perform consistently well or poorly on specific problems. In contrast, the aggregate distribution (left panel) is

not bimodal, indicating that different models perform well on different sets of problems. The availability of error feedback shifts the distribution of the sample-wise difficulty to the right, as expected. Interestingly, some samples are not solved at all across models, even with feedback, and no samples are solved consistently by all models. As a further investigation, we plan to qualitatively examine samples that shift from unsolved to solved given error feedback. Finally, the right panel shows that many samples are either "easy" or "hard", however larger models tend to have more "medium" difficulty samples, indicating that scale can, at least partially, improve version-conditioned generation from unsolved to solved. We qualitatively demonstrate some of these examples in A.1.

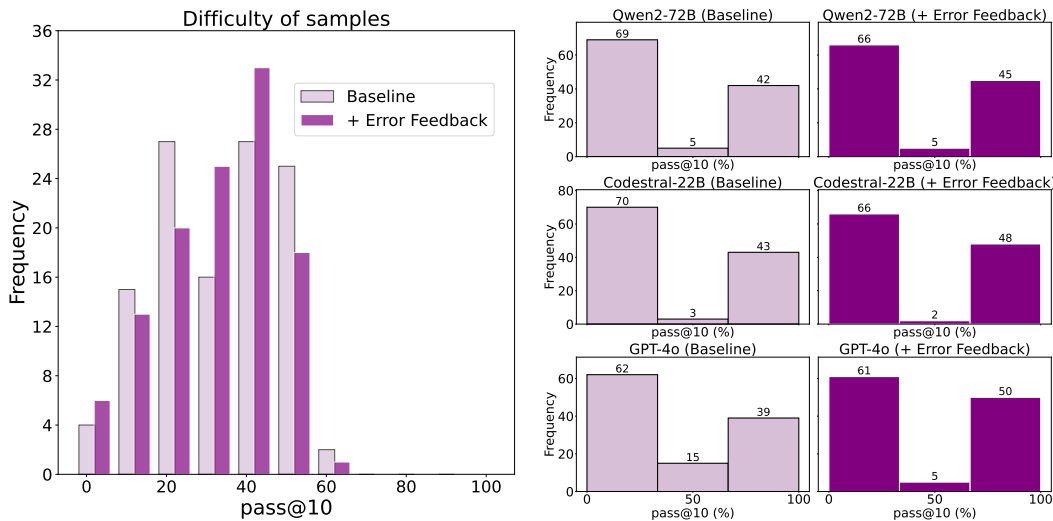


Figure 5: **Comparison of sample and model differences.** The left panel shows the distribution of sample difficulty, measured by the frequency of pass@10 scores across all models and seeds. The right panel presents the same distribution, but averaged for specific models across their seeds. Interestingly, individual models tend to exhibit bimodal distributions, indicating they are either consistently good or bad at specific problems. However, the aggregate distribution is not bimodal, suggesting that different models excel at different problems.

4 RELATED WORK

Code LLM training and evaluation protocols Code LLM evaluations mainly revolve around code completion (Zhang et al. (2023); van Dam et al. (2023); Lu et al. (2021)). Existing benchmarks emphasize generic code completion, yet a recognized limitation is the inability of code LLM to generate and complete code that requires library and project-level knowledge (Xu & Zhu, 2022), let alone version-level knowledge, which is vital for real-world software applications.

Recent initiatives address repository-level code understanding by LLMs (Bairi et al. (2023); Shrivastava et al. (2023a;b); Liu et al. (2023); Guo et al. (2024)). Attempts at library-level code generation (Zan et al. (2022)) and consideration of dependencies between files (Guo et al. (2024)) have been made. However, these efforts do not directly address the challenge of accommodating version-sensitive changes, adding complexity.

The core issue arises from models being trained on library code without explicit knowledge of library versions or their functional changes. Consequently, when tasked with generating code specifically compatible with a particular library version, there is a significant risk models often encounter failures.

Datasets Existing datasets like HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and MTPB (Nijkamp et al., 2022) provide sets of handwritten prompts and test cases to evaluate code generated by code LLM. However, these datasets are relatively small and lack context regarding a model’s comprehension of repositories. APPS (Hendrycks et al., 2021) and CodeContest (Li

et al., 2022) offer challenging datasets with coding competition questions, providing insights into a model’s performance on difficult problems but without a focus on library-specific challenges. DSP (Chandel et al., 2022) and DS-1000 (Lai et al., 2023) concentrate on the top data science libraries in Python, while JulCe (Agashe et al., 2019) uses Jupyter Notebooks for training and evaluation, but these notebooks do not necessarily need to be repository-specific. CoNaLa (Yin et al., 2018) contains problems collected from StackOverflow across multiple programming languages, including both library-specific questions and non-library-specific code. More recently, BigCodeBench (Zhuo et al., 2024) is constructed to evaluate the comprehensive capabilities of code generation with tool use and instruction following, which poses a great challenge for existing models. Several datasets include version-specific code, such as Versicode (Wu et al., 2024b), CodeUpdateArena (Liu et al., 2024), and the dataset by Wang et al. Versicode’s dataset, compiled from academic papers, Stack Overflow, and library source code, supports tasks like token, line, and block-level code completion and code editing. Unlike our dataset, Versicode evaluates using exact matches. Wang et al.’s dataset collects API mappings, such as “torch.lstsq() is deprecated in favor of torch.linalg.lstsq(),” and evaluates LLMs using exact match, edit similarity, and fixed rate metrics. Although Versicode and Wang et al.’s datasets address the evolving nature of libraries, their evaluations are limited to string matching.

In contrast, CodeUpdateArena evaluates LLMs’ ability to adapt to API changes, such as adding a boolean flag, by running tests. However, the dataset is synthetic and are not extracted from real-life version changes. For CodeUpdateArena, they also take the approach of training LLMs using the updated API function –using docstrings or examples–. It then tests if without access to the update during inference, the LLM’s reflects the synthetic changes. While these datasets provide valuable resources for training and evaluating models, our **GitChameleon** dataset advances research into version-conditioned code generation by LLMs. Runnable tests offer insights into LLM adaptability, and **GitChameleon** further assesses a model’s ability to differentiate between library versions, and successfully use a specific version.

Implications for Lifelong Learning Continual/lifelong learning in code generation models is in its early stages (Yadav et al., 2023; Weyssow et al., 2023; Wu et al., 2024a; Gao et al., 2023). However, current efforts often focus on artificial sequential tasks rather than utilizing the natural distribution shift in the chronological evolution of code. Notably, continual learning mainly targets mitigating catastrophic forgetting and balancing forward- and backward-transfer on a data stream, which may not align optimally with coding environment demands. In coding environments, obsolete or legacy libraries may prompt selective forgetting of irrelevant knowledge, particularly at the library/package level. Previous work, such as Caccia et al. (2021) may serve as a foundation for developing continual learning in Large Language Models for code.

5 LIMITATIONS

We consider the lack of prompt optimization done for the instruct models as a considerable limitation of our analysis. Furthermore, the dataset consists of 116 problems, which is relatively small compared to other code benchmarks. Finally, we do not explore approaches such as RAG, chain of thought, or finetuning on a split of our benchmark to observe an upper bound of performance on this task. Future work could explore such approaches using our dataset.

6 CONCLUSION

Recognizing the crucial need for code LLM adaptation to evolving code environments, particularly in widely used libraries, we introduce a novel and extensive Python-based version-specific benchmark named **GitChameleon**. By effectively leveraging **GitChameleon**, we expose the shortcomings of existing state-of-the-art (SoTA) models in producing version-specific code, representing an inaugural effort to draw attention to this challenge. While our work exposes this shortcoming, we acknowledge the dataset’s limitations. In future endeavors, we aim to enhance the dataset’s comprehensiveness across various programming languages and frameworks. Additionally, we plan to introduce new tasks that can benefit research on code LLM models using **GitChameleon**.

REFERENCES

- 540
541
542 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen
543 Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko,
544 Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dong-
545 dong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang
546 Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit
547 Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao,
548 Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin
549 Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim,
550 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden,
551 Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong
552 Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro
553 Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-
554 Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo
555 de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim,
556 Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla,
557 Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua
558 Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp
559 Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Ji-
560 long Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan,
561 Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan
562 Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your
563 phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 564
565 Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou.
566 Gradio: Hassle-free sharing and testing of ml models in the wild, 2019. URL <https://arxiv.org/abs/1906.02569>.
- 567
568 Rajas Agashe, Srinivasan Iyer, and Luke Zettlemoyer. Juice: A large scale distantly supervised
569 dataset for open domain context-based code generation. *arXiv preprint arXiv:1910.02216*, 2019.
- 570
571 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
572 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
573 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 574
575 Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Vageesh D. C, Arun Iyer, Suresh
576 Parthasarathy, Sriram Rajamani, B. Ashok, and Shashank Shet. CodePlan: Repository-level Cod-
577 ing using LLMs and Planning, September 2023. URL [http://arxiv.org/abs/2309.](http://arxiv.org/abs/2309.12499)
578 12499. arXiv:2309.12499 [cs].
- 579
580 Steven Bird and Edward Loper. NLTK: The natural language toolkit. In *Proceedings of the ACL*
581 *Interactive Poster and Demonstration Sessions*, pp. 214–217, Barcelona, Spain, July 2004. Asso-
582 ciation for Computational Linguistics. URL <https://aclanthology.org/P04-3031>.
- 583
584 Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel,
585 Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake Van-
586 derPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. API design for machine learning software:
587 experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Min-*
588 *ing and Machine Learning*, pp. 108–122, 2013.
- 589
590 Massimo Caccia, Pau Rodriguez, Oleksiy Ostapenko, Fabrice Normandin, Min Lin, Lucas Cac-
591 cia, Issam Laradji, Irina Rish, Alexandre Lacoste, David Vazquez, and Laurent Charlin. Online
592 Fast Adaptation and Knowledge Accumulation: a New Approach to Continual Learning, January
593 2021. URL <http://arxiv.org/abs/2003.05856>. arXiv:2003.05856 [cs].
- 594
595 Shubham Chandel, Colin B Clement, Guillermo Serrato, and Neel Sundaresan. Training and evalu-
596 ating a jupyter notebook data science assistant. *arXiv preprint arXiv:2201.12901*, 2022.
- 597
598 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
599 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
600 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,

- 594 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
595 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-
596 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex
597 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
598 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec
599 Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-
600 Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large
601 language models trained on code. 2021.
- 602 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and
603 memory-efficient exact attention with io-awareness, 2022. URL [https://arxiv.org/abs/
604 2205.14135](https://arxiv.org/abs/2205.14135).
- 605 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
606 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
607 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
608 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
609 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
610 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
611 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
612 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
613 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
614 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
615 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
616 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
617 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
618 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
619 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
620 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
621 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
622 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der
623 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
624 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
625 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
626 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
627 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
628 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
629 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
630 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
631 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
632 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
633 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
634 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
635 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
636 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
637 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
638 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
639 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
640 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
641 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
642 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
643 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
644 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
645 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
646 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
647 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,

- 648 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
649 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
650 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
651 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
652 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
653 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
654 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
655 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
656 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
657 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
658 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
659 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
660 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
661 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
662 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
663 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
664 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
665 Khabisa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
666 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
667 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
668 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
669 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
670 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
671 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
672 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
673 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
674 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
675 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
676 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
677 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
678 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
679 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
680 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
681 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
682 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
683 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
684 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
685 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
686 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
687 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
688 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
689 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.
690 URL <https://arxiv.org/abs/2407.21783>.
- 688 Shuzheng Gao, Hongyu Zhang, Cuiyun Gao, and Chaozheng Wang. Keeping pace with ever-
689 increasing data: Towards continual learning of code intelligence models, 2023.
690
- 691 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
692 Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the
693 large language model meets programming – the rise of code intelligence, 2024. URL <https://arxiv.org/abs/2401.14196>.
694
- 695 Aric Hagberg, Pieter J Swart, and Daniel A Schult. Exploring network structure, dynamics, and
696 function using networkx. Technical report, Los Alamos National Laboratory (LANL), Los
697 Alamos, NM (United States), 2008.
698
- 699 Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen,
700 David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert
701 Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane,
Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard,

- 702 Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Ar-
703 ray programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/
704 s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
705
- 706 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
707 Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence
708 with apps. *arXiv preprint arXiv:2105.09938*, 2021.
- 709 J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):
710 90–95, 2007. doi: 10.1109/MCSE.2007.55.
711
- 712 Kelsey Jordahl, Joris Van den Bossche, Martin Fleischmann, Jacob Wasserman, James McBride,
713 Jeffrey Gerard, Jeff Tratner, Matthew Perry, Adrian Garcia Badaracco, Carson Farmer, Geir Arne
714 Hjelle, Alan D. Snow, Micah Cochran, Sean Gillies, Lucas Culbertson, Matt Bartos, Nick Eu-
715 bank, maxalbert, Aleksey Bilogur, Sergio Rey, Christopher Ren, Dani Arribas-Bel, Leah Wasser,
716 Levi John Wolf, Martin Journois, Joshua Wilson, Adam Greenhall, Chris Holdgraf, Filipe, and
717 François Leblanc. *geopandas/geopandas: v1.0.1*, July 2020. URL <https://doi.org/10.5281/zenodo.3946761>.
718
- 719 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau
720 Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data
721 science code generation. In *International Conference on Machine Learning*, pp. 18319–18345.
722 PMLR, 2023.
- 723 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao
724 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii,
725 Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João
726 Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Lo-
727 gesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra
728 Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey,
729 Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luc-
730 cioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor,
731 Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex
732 Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva
733 Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes,
734 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. *StarCoder: may the source
735 be with you!*, 2023.
- 736 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
737 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien
738 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven
739 Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson,
740 Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level
741 code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022. ISSN 1095-
742 9203. doi: 10.1126/science.abq1158. URL <http://dx.doi.org/10.1126/science.abq1158>.
743
- 744 Tianyang Liu, Canwen Xu, and Julian McAuley. *RepoBench: Benchmarking Repository-Level
745 Code Auto-Completion Systems*, October 2023. URL <http://arxiv.org/abs/2306.03091>. arXiv:2306.03091 [cs].
746
- 747 Zeyu Leo Liu, Shrey Pandit, Xi Ye, Eunsol Choi, and Greg Durrett. *Codeupdatearena: Bench-
748 marking knowledge editing on api updates*, 2024. URL <https://arxiv.org/abs/2407.06249>.
749
- 750 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
751 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. *StarCoder 2 and the stack v2: The
752 next generation*. *arXiv preprint arXiv:2402.19173*, 2024.
753
- 754 Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin
755 Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou,
Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu

- 756 Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code understanding
757 and generation, 2021.
758
- 759 Meta. Llama 3: Meta’s latest large language model. [https://github.com/meta-llama/
760 llama3](https://github.com/meta-llama/llama3), 2024. Accessed: 2024-06-03.
761
- 762 Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza So-
763 ria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, Manish Sethi,
764 Xuan-Hong Dang, Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew Coleman, Matthew White,
765 Mark Lewis, Raju Pavuluri, Yan Koyfman, Boris Lublinsky, Maximilien de Baysier, Ibrahim
766 Abdelaziz, Kinjal Basu, Mayank Agarwal, Yi Zhou, Chris Johnson, Aanchal Goyal, Hima Pa-
767 tel, Yousaf Shah, Petros Zerfos, Heiko Ludwig, Asim Munawar, Maxwell Crouse, Pavan Ka-
768 panipathi, Shweta Salaria, Bob Calio, Sophia Wen, Seetharami Seelam, Brian Belgodere, Car-
769 los Fonseca, Amith Singhee, Nirmal Desai, David D. Cox, Ruchir Puri, and Rameswar Panda.
770 Granite code models: A family of open foundation models for code intelligence, 2024. URL
771 <https://arxiv.org/abs/2405.04324>.
772
- 773 MistralAI. Codestral: Hello, world! 2024.
774
- 775 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and
776 Caiming Xiong. A conversational paradigm for program synthesis. *arXiv preprint*, 2022.
777
- 778 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
779 cia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red
780 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
781 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher
782 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
783 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann,
784 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
785 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
786 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
787 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila
788 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
789 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-
790 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan
791 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-
792 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan
793 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu,
794 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
795 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
796 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook
797 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel
798 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
799 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel
800 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez,
801 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv
802 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney,
803 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,
804 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel
805 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-
806 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
807 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
808 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
809 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny,
Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,
Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra
Rimbach, Carl Ross, Bob Rotsted, Henri Roussee, Nick Ryder, Mario Saltarelli, Ted Sanders,
Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-
sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor,
Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,

- 810 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
811 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-
812 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-
813 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan
814 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lillian Weng,
815 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-
816 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
817 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
818 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
819 <https://arxiv.org/abs/2303.08774>.
- 820 The pandas development team. pandas-dev/pandas: Pandas, February 2020. URL [https://doi.](https://doi.org/10.5281/zenodo.3509134)
821 [org/10.5281/zenodo.3509134](https://doi.org/10.5281/zenodo.3509134).
- 822 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
823 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-
824 ward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
825 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep
826 learning library, 2019.
- 827 Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu,
828 and Nathan Cooper. Stable code 3b. URL [[https://huggingface.co/](https://huggingface.co/stabilityai/stable-code-3b)
829 [stabilityai/stable-code-3b](https://huggingface.co/stabilityai/stable-code-3b)] ([https://huggingface.co/stabilityai/](https://huggingface.co/stabilityai/stable-code-3b)
830 [stable-code-3b](https://huggingface.co/stabilityai/stable-code-3b)).
- 831 Alibaba Qwen. Qwen2 technical report. 2024.
- 832 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
833 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton,
834 Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez,
835 Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and
836 Gabriel Synnaeve. Code Llama: Open Foundation Models for Code, August 2023. URL [http:](http://arxiv.org/abs/2308.12950)
837 arxiv.org/abs/2308.12950. arXiv:2308.12950 [cs].
- 838 Disha Shrivastava, Denis Kocetkov, Harm de Vries, Dzmitry Bahdanau, and Torsten Scholak. Re-
839 poFusion: Training Code Models to Understand Your Repository, June 2023a. URL [http:](http://arxiv.org/abs/2306.10998)
840 arxiv.org/abs/2306.10998. arXiv:2306.10998 [cs].
- 841 Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. Repository-Level Prompt Generation
842 for Large Language Models of Code, June 2023b. URL [http://arxiv.org/abs/2206.](http://arxiv.org/abs/2206.12839)
843 [12839](http://arxiv.org/abs/2206.12839). arXiv:2206.12839 [cs].
- 844 CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu,
845 Christopher A. Choquette-Choo, Jingyue Shen, Joe Kelley, Kshitij Bansal, Luke Vilnis, Mateo
846 Wirth, Paul Michel, Peter Choy, Pratik Joshi, Ravin Kumar, Sarmad Hashmi, Shubham Agrawal,
847 Zhitao Gong, Jane Fine, Tris Warkentin, Ale Jakse Hartman, Bin Ni, Kathy Korevec, Kelly
848 Schaefer, and Scott Huffman. Codegemma: Open code models based on gemma, 2024. URL
849 <https://arxiv.org/abs/2406.11409>.
- 850 Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of con-
851 text, 2024a. URL <https://arxiv.org/abs/2403.05530>.
- 852 Qwen Team. Code with codeqwen1.5, April 2024b. URL [https://qwenlm.github.io/](https://qwenlm.github.io/blog/codeqwen1.5/)
853 [blog/codeqwen1.5/](https://qwenlm.github.io/blog/codeqwen1.5/).
- 854 Qwen Team. Qwen2.5: A party of foundation models, September 2024c. URL [https:](https://qwenlm.github.io/blog/qwen2.5/)
855 qwenlm.github.io/blog/qwen2.5/.
- 856 Tim van Dam, Maliheh Izadi, and Arie van Deursen. Enriching source code with contextual data for
857 code completion models: An empirical study, 2023.
- 858 Chong Wang, Kaifeng Huang, Jian Zhang, Yebo Feng, Lyuye Zhang, Yang Liu, and Xin Peng.
859 How and why llms use deprecated apis in code completion? an empirical study, 2024. URL
860 <https://arxiv.org/abs/2406.09834>.

- 864 Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and
865 Jarrod Millman (eds.), *Proceedings of the 9th Python in Science Conference*, pp. 56 – 61, 2010.
866 doi: 10.25080/Majora-92bf1922-00a.
867
- 868 Martin Weyssow, Xin Zhou, Kisub Kim, David Lo, and Houari Sahraoui. On the usage of continual
869 learning for out-of-distribution generalization in pre-trained language models of code. In *Pro-*
870 *ceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on*
871 *the Foundations of Software Engineering*, pp. 1470–1482, 2023.
- 872 Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari.
873 Continual learning for large language models: A survey, 2024a.
- 874 Tongtong Wu, Weigang Wu, Xingyu Wang, Kang Xu, Suyu Ma, Bo Jiang, Ping Yang, Zhenchang
875 Xing, Yuan-Fang Li, and Gholamreza Haffari. Versicode: Towards version-controllable code
876 generation, 2024b. URL <https://arxiv.org/abs/2406.07411>.
- 877 Yichen Xu and Yanqiao Zhu. A Survey on Pretrained Language Models for Neural Code Intelli-
878 gence, December 2022. URL <http://arxiv.org/abs/2212.10079>. arXiv:2212.10079
879 [cs].
- 881 Prateek Yadav, Qing Sun, Hantian Ding, Xiaopeng Li, Dejiao Zhang, Ming Tan, Xiaofei Ma, Par-
882 minder Bhatia, Ramesh Nallapati, Murali Krishna Ramanathan, Mohit Bansal, and Bing Xi-
883 ang. Exploring Continual Learning for Code Generation Models, July 2023. URL [http://](http://arxiv.org/abs/2307.02435)
884 arxiv.org/abs/2307.02435. arXiv:2307.02435 [cs].
885
- 886 Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. Learning to
887 mine aligned code and natural language pairs from stack overflow. In *Proceedings of the 15th*
888 *international conference on mining software repositories*, pp. 476–486, 2018.
- 889 Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang,
890 Weizhu Chen, and Jian-Guang Lou. CERT: Continual Pre-Training on Sketches for Library-
891 Oriented Code Generation, June 2022. URL <http://arxiv.org/abs/2206.06888>.
892 arXiv:2206.06888 [cs].
- 893 Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang
894 Lou, and Weizhu Chen. RepoCoder: Repository-Level Code Completion Through Iterative
895 Retrieval and Generation, October 2023. URL <http://arxiv.org/abs/2303.12570>.
896 arXiv:2303.12570 [cs].
897
- 898 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari,
899 Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Bench-
900 marking code generation with diverse function calls and complex instructions. *arXiv preprint*
901 *arXiv:2406.15877*, 2024.

903 A APPENDIX

904 A.1 QUALITATIVE SAMPLES OF ERROR FEEDBACK

905 Listing 1: Original output from Codestral-22B-v0.1

```

906
907
908
909 1 # [problem]
910 2 Calculate the natural logarithm of the absolute value of the gamma
    function using pytorch special functions if available in this version
911  , otherwise you may use another library.
912 3 # [starting code]
913 4 import torch
914 5 input_tensor = torch.linspace(0, 10, steps=10)
915 6 output =
916 7 # [sample solution]
917 8 torch.special.gammaln(input_tensor)
918 9
919 10 # [model output]
```

```

918 11 import torch
919 12 input_tensor = torch.linspace(0, 10, steps=10)
920 13 output = torch.log(torch.abs(torch.special.gamma(input_tensor)))
921 14 # [Error trace]
922 15 Traceback (most recent call last):
923 16     output = torch.log(torch.abs(torch.special.gamma(input_tensor)))
924 17 AttributeError: module 'torch.special' has no attribute 'gamma'. Did you
    mean: 'digamma'?

```

Listing 2: Regenerated solution with Codestral-22B-v0.1 given error feedback

```

928 1 import torch
929 2 def natural_log_abs_gamma(input_tensor):
930 3     gamma_tensor = torch.lgamma(input_tensor)
931 4     natural_log_abs_gamma_tensor = torch.abs(gamma_tensor)
932 5     return natural_log_abs_gamma_tensor
933 6 input_tensor = torch.linspace(0, 10, steps=10)
934 7 output = natural_log_abs_gamma(input_tensor)

```

Listing 3: Original output from Phi-3.5-MoE-instruct

```

935 1 # torch-1.10.0
936 2 # [problem]
937 3 Calculate the error function using pytorch special functions if available
938 4 in this version, otherwise you may use another library.
939 5 # [starting code]
940 6 import torch\ninput_tensor = torch.linspace(0, 10, steps=10)\noutput =
941 7 # [sample solution]
942 8 torch.special.erf(input_tensor)
943 9
944 10 # [model output]
945 11 import torch
946 12 from scipy.special import erf
947 13 input_tensor = torch.linspace(0, 10, steps=10)
948 14 def calculate_erf(tensor):
949 15     # Convert tensor to numpy array
950 16     numpy_array = tensor.numpy()
951 17     # Calculate erf using scipy.special.erf
952 18     output = erf(numpy_array)
953 19     # Convert numpy array back to tensor
954 20     output_tensor = torch.from_numpy(output)
955 21     return output_tensor
956 22 output = calculate_erf(input_tensor)
957 23 # [Error trace]
958 24 Traceback (most recent call last):
959 25     from scipy.special import erf
960 ModuleNotFoundError: No module named 'scipy'

```

Listing 4: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```

961 1 import torch
962 2 def compute_error_function(input_tensor):
963 3     output = torch.erf(input_tensor)
964 4     return output
965 5 input_tensor = torch.linspace(0, 10, steps=10)
966 6 output = compute_error_function(input_tensor)

```

Listing 5: Original output with Phi-3.5-MoE-instruct

```

969 1 # numpy-1.21.0
970 2 # [problem]
971 3 Write a function that rounds an array of numbers.
972 4 # [starting code]

```

```

972 5 import numpy as np
973 6
974 7 def custom_round(arr):
975 8     return
976 9 # [sample solution]
977 10 np.round_(arr)
978 11 # [model output]
979 12
979 13 import numpy as np
980 14 def custom_round(arr):
981 15     return np.round(arr)
982 16 rounded_arr = custom_round(arr)
983 17
984 18
984 19 def test_custom_round():
985 20     arr = np.array([1.5, 2.3, 3.7])
986 21     result = custom_round(arr)
987 22     expected = np.round_(arr)
988 23     assert np.array_equal(result, expected)
989 24
989 25 test_custom_round()
990 26 # [Error trace]
991 27 Traceback (most recent call last):
992 28     rounded_arr = custom_round(arr)
993 29 NameError: name 'arr' is not defined

```

Listing 6: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```

996 1 import numpy as np
997 2 def custom_round(arr):
998 3     return np.round(arr)
999

```

Listing 7: Original output with Phi-3.5-MoE-instruct

```

1002 1 # numpy-1.25.0
1003 2 # [problem]
1004 3 Write a function that rounds an array of numbers.
1005 4 # [starting code]
1006 5 import numpy as np
1007 6
1007 7 def custom_round(arr):
1008 8     return
1009 9 # [sample solution]
1010 10 np.round(arr)
1011 11 # [model output]
1012 12
1012 13 import numpy as np
1013 14 def custom_round(arr):
1014 15     return np.round_(arr)
1015 16
1016 17
1016 18 def test_custom_round():
1017 19     arr = np.array([1.5, 2.3, 3.7])
1018 20     result = custom_round(arr)
1019 21     expected = np.round(arr)
1020 22     assert np.array_equal(result, expected)
1021 23
1021 24 test_custom_round()
1022 25 # [Error trace]
1023 26 DeprecationWarning: `round_` is deprecated as of NumPy 1.25.0, and will
1024 27     be removed in NumPy 2.0. Please use `round` instead.
1025 27     result = custom_round(arr)

```

1026 Listing 8: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```
1027
1028 1 import numpy as np
1029 2 def custom_round(arr):
1030 3     return np.round(arr)
```

1031 Listing 9: Original output with Phi-3.5-MoE-instruct

```
1033 1 # pandas-1.4.0
1034 2 # [problem]
1035 3 Predict the correct type.
1036 4 # [starting code]
1037 5 import pandas as pd\nindex = pd.Index([1, 2, 3], dtype='int32')\
1038 6     nis_correct_type = index.dtype ==
1039 7     # [sample solution]
1040 8     'int64'
1041 9 # [model output]
1042 10 import pandas as pd
1043 11 def check_dtype(index):
1044 12     is_correct_type = check_dtype(index)
1045 13 # [Error trace]
1046 14 Sorry: IndentationError: expected an indented block after function
1047     definition on line 3
```

1048 Listing 10: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```
1049 1 import pandas as pd
1050 2 def is_correct_type(index, expected_dtype):
1051 3     return index.dtype == expected_dtype
1052 4 index = pd.Index([1, 2, 3], dtype='int32')
```

1053 Listing 11: Original output with Phi-3.5-MoE-instruct

```
1055 1 # gradio-2.9.2
1056 2 # [problem]
1057 3 Write a function that takes an image input and returns a textbox output.
1058 4 # [starting code]
1059 5 import gradio as gr
1060 6 def process_image(image):
1061 7     return "Processed"
1062 8
1063 9 iface = gr.Interface
1064 10 # [sample solution]
1065 11 (fn=process_image, inputs=gr.inputs.Image(), outputs=gr.outputs.Textbox()
1066 12     )
1067 13 # [model output]
1068 14 import gradio as gr
1069 15 def process_image(image):
1070 16     # Perform processing on the image here
1071 17     # For simplicity, we're just returning the image dimensions
1072 18     width, height = image.shape[1], image.shape[0]
1073 19     return f"Width: {width}, Height: {height}"
1074 20 iface = gr.Interface(
1075 21     fn=process_image,
1076 22     inputs=gr.Image(),
1077 23     outputs="text"
1078 24 )
1079 25 # [Error trace]
1080 26 Traceback (most recent call last):
1081 27     inputs=gr.Image(),
1082 28     AttributeError: module 'gradio' has no attribute 'Image'
```


Listing 12: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```

1080
1081
1082 1 import gradio as gr
1083 2 from PIL import Image
1084 3 import numpy as np
1085 4 def process_image(image):
1086 5     # Convert gradio image to PIL Image
1087 6     pil_image = Image.open(image)
1088 7
1089 8     # Process the image (example: convert to grayscale)
1090 9     processed_image = pil_image.convert('L')
1091 10
1092 11     # Convert processed image back to numpy array
1093 12     processed_image_array = np.array(processed_image)
1094 13
1095 14     # Return the processed image array as a string
1096 15     return str(processed_image_array.tolist())
1097 16
1098 17 iface = gr.Interface(
1099 18     fn=process_image,
1100 19     inputs=gr.inputs.Image(shape=(None, None, 3)),
1101 20     outputs=gr.outputs.Textbox()
1102 )

```

Listing 13: Original output with Phi-3.5-MoE-instruct

```

1101 1 # networkx-2.5
1102 2 # [problem]
1103 3
1104 4 Write a function that returns the naive greedy modularity communities for
1105 5 a graph.
1106 6
1107 7 # [starting code]
1108 8 import networkx as nx
1109 9 def naive_modularity_communities(G):
1110 10     return nx.community.
1111 11 # [sample solution]
1112 12 naive_greedy_modularity_communities(G)
1113 13 # [model output]
1114 14
1115 15 import networkx as nx
1116 16 def naive_modularity_communities(G):
1117 17     communities = []
1118 18     while G.number_of_nodes() > 0:
1119 19         node = max(G.degree, key=lambda x: x[1])[0]
1120 20         community = nx.algorithms.community.greedy_modularity_communities
1121 21             (G, nbunch=[node])
1122 22         communities.append(community)
1123 23         G.remove_nodes_from(community)
1124 24     return communities
1125 25
1126 26 G = nx.karate_club_graph()
1127 27 communities = naive_modularity_communities(G)
1128 28
1129 29 G = nx.karate_club_graph()
1130 30 # [Error trace]
1131 31 Traceback (most recent call last):
1132 32     community = nx.algorithms.community.greedy_modularity_communities(G,
1133 33         nbunch=[node])
1134 34 TypeError: greedy_modularity_communities() got an unexpected keyword
1135 35 argument 'nbunch'

```

Listing 14: Regenerated solution with Phi-3.5-MoE-instruct given error feedback

```

1132 1 import networkx as nx
1133 2 def naive_modularity_communities(G):
1134 3     communities = []

```

```

1134 4   unvisited_nodes = list(G.nodes)
1135 5
1136 6   while unvisited_nodes:
1137 7       current_node = unvisited_nodes[0]
1138 8       community = nx.algorithms.community.greedy_modularity_communities
1139 9           (G, start_nodes=[current_node])
1140 10      visited_nodes = list(community[0])
1141 11      unvisited_nodes = list(set(unvisited_nodes) - set(visited_nodes))
1142 12      communities.append(visited_nodes)
1143 13
1144 14      return communities
1145 15  G = nx.karate_club_graph()
1146      result = naive_modularity_communities(G)

```

Listing 15: Original output with gpt4o

```

1148 1 # torch-1.10.0
1149 2 # [problem]
1150 3 Calculate the complementary error function using pytorch's special
1151 4     functions if available in this version, otherwise you may use another
1152 5     library.
1153 6 # [starting code]
1154 7 import torch\ninput_tensor = torch.linspace(0, 10, steps=10)\noutput =
1155 8 # [sample solution]
1156 9 torch.special.erfc(input_tensor)
1157 10 # [model output]
1158 11
1159 12 import torch
1160 13 import scipy.special
1161 14 def complementary_error_function(x):
1162 15     # Check if 'special.erfc' is available in the current version of
1163 16     PyTorch
1164 17     if hasattr(torch.special, 'erfc'):
1165 18         return torch.special.erfc(x)
1166 19     else:
1167 20         # Use scipy's erfc function if PyTorch's is not available
1168 21         x_np = x.detach().numpy() # Convert the tensor to a numpy array
1169 22         erfc_np = scipy.special.erfc(x_np) # Calculate erfc using scipy
1170 23         return torch.from_numpy(erfc_np) # Convert the result back to a
1171 24         tensor
1172 25 input_tensor = torch.linspace(0, 10, steps=10)
1173 26 output = complementary_error_function(input_tensor)
1174 27 # [Error trace]
1175 28 Traceback (most recent call last):
1176 29   import scipy.special
1177 30 ModuleNotFoundError: No module named 'scipy'

```

Listing 16: Regenerated solution with gpt4o given error feedback

```

1177 1 import torch
1178 2 def complementary_error_function(x):
1179 3     # Check if 'special.erfc' is available in the current version of
1180 4     PyTorch
1181 5     if hasattr(torch.special, 'erfc'):
1182 6         return torch.special.erfc(x)
1183 7     else:
1184 8         # Raise an error if PyTorch's special.erfc is not available and
1185 9         scipy is not installed
1186 10        raise ImportError("PyTorch version does not support special.erfc
1187 11        and scipy is not installed.")
1188 12 input_tensor = torch.linspace(0, 10, steps=10)
1189 13 output = complementary_error_function(input_tensor)

```

Listing 17: Original output with gpt4o

```

1188
1189
1190 1 # scikit-learn-1.3
1191 2 # [problem]
1192 3 Retrieve and list all available scorer names, ensuring they are returned
1193 4 in a list format.
1194 5 # [starting code]
1195 6 from sklearn import metrics\nscorer_names =
1196 7 # [sample solution]
1197 8 metrics.get_scorer_names()
1198 9 # [model output]
1199 10
1200 11 from sklearn import metrics
1201 12 scorer_names = list(metrics.SCORERS.keys())
1202 13 # [Error trace]
1203 14 Traceback (most recent call last):
1204 15 AttributeError: module 'sklearn.metrics' has no attribute 'SCORERS'

```

Listing 18: Regenerated solution with gpt4o given error feedback

```

1205 1 from sklearn.metrics import get_scorer_names
1206 2 scorer_names = list(get_scorer_names())
1207 3 conditions = isinstance(scorer_names, list) and len(scorer_names) > 0

```

Listing 19: Original output with gpt4o

```

1210 1 # gradio-3.36.0
1211 2 # [problem]
1212 3 Write a function that displays an image using Gradio where you cannot
1213 4 share the image.
1214 5 # [starting code]
1215 6 import gradio as gr
1216 7 def display_image():
1217 8     return "https://image_placeholder.com/42"
1218 9
1219 10 iface = gr.Interface
1220 11 # [sample solution]
1221 12 (fn=display_image, inputs=[], outputs=gr.Image(show_share_button=False))
1222 13 # [model output]
1223 14
1224 15 import gradio as gr
1225 16 def display_image():
1226 17     # This is just a placeholder URL; replace it with your actual image
1227 18     # URL or file path.
1228 19     return "https://via.placeholder.com/150"
1229 20 iface = gr.Interface(
1230 21     fn=display_image,
1231 22     inputs=None,
1232 23     outputs=gr.Image(label="Display Image", type="auto"),
1233 24     live=False # Ensure live updates are off
1234 25 )
1235 26 # [Error trace]
1236 27 line 116, in __init__
1237 28 raise ValueError(
1238 29     ValueError: Invalid value for parameter 'type': auto. Please choose from
1239 30     one of: ['numpy', 'pil', 'filepath']

```

Listing 20: Regenerated solution with gpt4o given error feedback

```

1239 1 import gradio as gr
1240 2 def display_image():
1241 3     # This is just a placeholder URL; replace it with your actual image
1242 4     # URL or file path.

```

```

1242 4     return "https://via.placeholder.com/150"
1243 5 iface = gr.Interface(
1244 6     fn=display_image,
1245 7     inputs=None,
1246 8     outputs=gr.Image(label="Display Image", type="filepath"),
1247 9     live=False # Ensure live updates are off
1248 10 )
1249 11 iface.share = False

```

Listing 21: Original output with gpt4o

```

1251 1 # gradio-3.24.0
1252 2 # [problem]
1253 3 Write a function that takes an image input and returns a label output.
1254 4 # [starting code]
1255 5 import gradio as gr
1256 6 def process_image(image):
1257 7     return "Processed"
1258 8
1259 9 iface = gr.Interface
1260 10 # [sample solution]
1261 11 (fn=process_image, inputs=gr.Image(), outputs=gr.Label())
1262 12 # [model output]
1263 13
1264 14 import gradio as gr
1265 15 def process_image(image):
1266 16     # Placeholder function to process the image and return a label
1267 17     # You can add your image processing and model inference logic here
1268 18     return "Processed"
1269 19 iface = gr.Interface(
1270 20     fn=process_image, # The function to process the image
1271 21     inputs=gr.inputs.Image(), # The input component (an image)
1272 22     outputs=gr.outputs.Text() # The output component (a text label)
1273 23 )
1274 24 # [Error trace]
1275 25 python3.10/site-packages/gradio/inputs.py:257: UserWarning: Usage of
1276 26 gradio.inputs is deprecated, and will not be supported in the future,
1277 27 please import your component from gradio.components
1278 28 warnings.warn(
1279 29 python3.10/site-packages/gradio/deprecation.py:40: UserWarning: 'optional
1280 30 ' parameter is deprecated, and it has no effect
1281 31 warnings.warn(value)
1282 32 Traceback (most recent call last):
1283 33   outputs=gr.outputs.Text() # The output component (a text label)
1284 34 AttributeError: module 'gradio.outputs' has no attribute 'Text'

```

Listing 22: Regenerated solution with gpt4o given error feedback

```

1283 1 import gradio as gr
1284 2 def process_image(image):
1285 3     # Placeholder function to process the image and return a label
1286 4     # You can add your image processing and model inference logic here
1287 5     return "Processed"
1288 6 iface = gr.Interface(
1289 7     fn=process_image, # The function to process the image
1290 8     inputs=gr.Image(), # The input component (an image)
1291 9     outputs=gr.Label() # The output component (a text label)
1292 10 )

```

Listing 23: Original output with gpt4o

```

1294 1 # networkx-2.5
1295 2 # [problem]
1296 3

```

```

1296 4 Write a function that accesses the first edge of a NetworkX graph.
1297 5
1298 6 # [starting code]
1299 7 import networkx as nx
1300 8 def get_first_edge(G):
1301 9     return
1302 10 # [sample solution]
1303 11 list(G.edges)[0]
1304 12 # [model output]
1305 13
1306 14 import networkx as nx
1307 15 def get_first_edge(G):
1308 16     # Check if the graph has any edges
1309 17     if G.number_of_edges() == 0:
1310 18         return None
1311 19     # Access the first edge
1312 20     first_edge = next(G.edges(data=True))
1313 21     return first_edge
1314 22
1315 23 G = nx.karate_club_graph()
1316 24 # [Error trace]
1317 25 Traceback (most recent call last):
1318 26   line 8, in get_first_edge
1319 27     first_edge = next(G.edges(data=True))
1320 28   TypeError: 'EdgeDataView' object is not an iterator
1321

```

Listing 24: Regenerated solution with gpt4o given error feedback

```

1320 1 import networkx as nx
1321 2 def get_first_edge(G):
1322 3     # Check if the graph has any edges
1323 4     if G.number_of_edges() == 0:
1324 5         return None
1325 6     # Access the first edge
1326 7     first_edge = list(G.edges(data=True))[0]
1327 8     return first_edge
1328 9 G = nx.karate_club_graph()
1329 10 first_edge = get_first_edge(G)

```

A.2 PSEUDOCODE FOR OUR VERIFICATION PROCESS

```

1332 1 # Create and activate virtual environment
1333 2 Run: "python -m venv venv"
1334 3 Run: "source venv/bin/activate"
1335 4
1336 5 # Install specified library and version
1337 6 Run: "pip install $library==$version"
1338 7
1339 8 # Install additional dependencies if specified
1340 9 If additional_dep:
1341 10     Run: "pip install $additional_dep"
1342 11
1343 12 # Combine code snippets
1344 13 complete_code = starter_code + expected_output + test
1345 14
1346 15 # Run the combined code with a timeout
1347 16 Run: "timeout 60 python -c '$complete_code'"
1348 17
1349 18 # Capture and print exit code
1350 19 exit_code = LastCommandExitCode()
1351 20 Print: "THIS WAS THE EXIT CODE: $exit_code"
1352 21
1353 22 # Print the complete code

```

```

1350 Print: complete_code
1351
1352 # Deactivate and remove virtual environment
1353 Run: "deactivate"
1354 Run: "rm -rf venv"

```

Each sample was validated using this method to ensure that it functioned as intended.

A.3 COMPARISON OF CODE LLMs

Table 5: Comparison of Code LLMs.

Model	Org.	Train. Cutoff Date	Pub. Avail. dataset
StarCoder (Li et al., 2023)	BigCode (HuggingFace, ServiceNow, NVIDIA)	04/2023	✓
StarCoder2 (Lozhkov et al., 2024)	BigCode (HuggingFace, ServiceNow, NVIDIA)	09/2023	✓
Qwen 2 (Qwen, 2024)	Alibaba	✗	✗
Qwen 2.5 (Team, 2024e)	Alibaba	✗	✗
CodeQwen 1.5 (Team, 2024b)	Alibaba	✗	✗
Codestral (MistralAI, 2024)	MistralAI	✗	✗
LLAMA3 (Meta, 2024)	Meta	03/2023, 12/2023	✗
LLAMA3.1 (Dubey et al., 2024)	Meta	✗	✗
LLAMA3.2	Meta	✗	✗
CodeLLAMA (Rozière et al., 2023)	Meta	10/2023	✗
DeepSeek-Coder (Guo et al., 2024)	DeepSeek-AI	02/2023	✗
CodeGemma (Team et al., 2024)	Google	✗	✗
Stable-Code (Pinnaparaju et al.)	Stability-AI	✗	✗
Granite-Code (Mishra et al., 2024)	IBM	✗	✗
Phi 3.5 (Abdin et al., 2024)	Microsoft	✗	✗
Nxcode-CQ	NTQA Solution	✗	✗
GPT-3.5	OpenAI	✗	✗
GPT-4 (OpenAI et al., 2024)	OpenAI	✗	✗
GPT-o1	OpenAI	✗	✗
Gemini 1.5 (Team, 2024a)	Google	✗	✗