

# EXECUTION-EVAL: CAN LANGUAGE MODELS EXECUTE REAL-WORLD CODE?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

As language models (LLMs) advance, traditional benchmarks face challenges of dataset saturation and disconnection from real-world performance, limiting our understanding of true model capabilities. We introduce EXecution-Eval (EXE), a benchmark designed to assess LLMs’ ability to execute code and predict program states. EXE attempts to address key limitations in existing evaluations: difficulty scaling, task diversity, training data contamination, and cost-effective scalability. Comprising over 30,000 tasks derived from 1,000 popular Python repositories on GitHub, EXE spans a wide range of lengths and algorithmic complexities. Tasks require models to execute code, necessitating various operations including mathematical reasoning, logical inference, bit manipulation, string operations, loop execution, and maintaining multiple internal variable states during computation. Our methodology involves: (a) selecting and preprocessing GitHub repositories, (b) generating diverse inputs for functions, (c) executing code to obtain ground truth outputs, and (d) formulating tasks that require models to reason about code execution. This approach allows for continuous new task generation for as few as 1,123 tokens, significantly reducing the risk of models “training on the test set.” We evaluate several state-of-the-art LLMs on EXE, revealing insights into their code comprehension and execution capabilities. Our results show that even the best-performing models struggle with complex, multi-step execution tasks, highlighting specific computational concepts that pose the greatest challenges for today’s LLMs. Furthermore, we review EXE’s potential for finding and predicting errors to aid in assessing a model’s cybersecurity capabilities. We propose EXE as a sustainable and challenging testbed for evaluating frontier models, offering insights into their internal mechanistic advancement.

## 1 INTRODUCTION

Language model benchmarks are facing challenges of rapid saturation (Ott et al., 2022) and an increasing disconnect from real-world performance perceived by end-users (Zheng et al., 2023). Due to this, benchmarks are being continually created to address failure modes; e.g. SuperGLUE targeting GLUE’s low problem difficulty (Wang et al., 2019), BIG-bench targeting general low eval diversity (Srivastava et al., 2022) and Auto-Arena-Hard targeting training-set contamination and data diversity in Chatbot-Arena (Li et al., 2024)(Chiang et al., 2024). These failure modes all demonstrate the challenge in linking the mechanistic improvements within language models to human understandable tasks.

Hence, to maximise an eval’s utility we aim to minimise the common failure modes of; a) difficulty, not ensuring an unbound scale of small trivial problems to complex multi-step problems, b) diversity, not ensuring a representative distribution across a large space of problems, c) novelty, not ensuring continually fresh, out-of-training data samples can be generated and, d) scalability, not ensuring tasks are cost-effective to generate in the thousands and beyond.

Motivated by these challenges we introduce EXecutionEval (EXE), an evaluation replicating one of the primary tasks humans perform while coding; predicting and comparing a final program state for a given set of inputs - seen in Figure 1. EXE is designed to avoid the aforementioned failure modes; emphasising difficulty (smooth scale from trivial 1-step, one-line functions to difficult 100s-of-step, multi-layer functions), diversity (unbound number of test cases generatable for tasks from

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

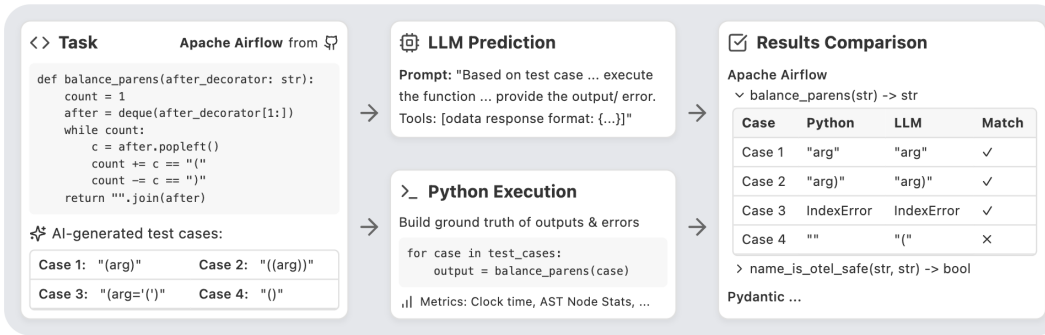


Figure 1: An example task from Apache Airflow’s Github repository (code simplified to fit within diagram). EXE sources tasks from 1,000 Python repositories, generates test cases for them, and compares the LLM’s ability to execute code against python’s interpreter.

1,000 GitHub Repos), novelty (program inputs can be continually generated) and scalability (initial release containing 30,000+ problems at a cost of \$11).

EXE also holds theoretical inspiration. (Fowler et al., 2022) et al have replicated positive pedagogical correlations found by (Lopez et al., 2008) between the abilities of CS1 students to “trace” programs (i.e. manually predict outputs and write the internal state out line by line) and their abilities to pass code writing and explanation exams. This is mirrored in CRUX-Eval’s (Gu et al., 2024) findings, where they observe a moderate correlation between a model’s ability to execute a block of code and a model’s HumanEval (Chen et al., 2021) code writing Pass@1 rate.

## 2 EVALUATION FRAMEWORK

As seen in Figure 1, an EXE task is to predict a function’s return value or error from: a) a code snippet and b) a set of input arguments. Code snippets are extracted from PyPi’s most popular 1,000 python projects hosted on GitHub, we select our snippets to be pure (i.e. deterministic, no side effects), language model generatable (i.e. arg types of ints, lists, ...) and to only require builtins (local imports and external libraries are inlined for the snippet). To realise this we follow the following three stage pipeline:

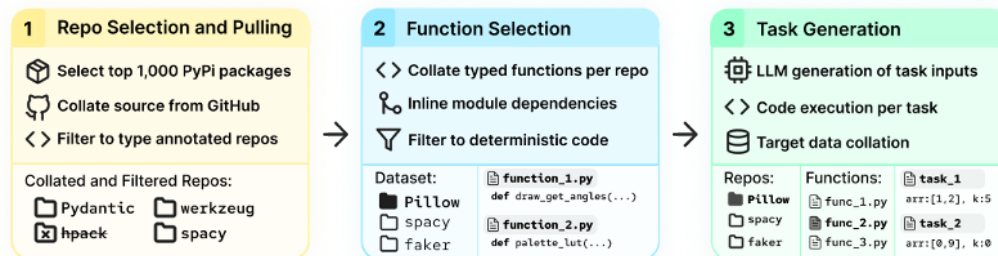


Figure 2: Three stage EXE task generation pipeline. Detailed example tasks and generated inputs can be found in Appendix A.1.

**1. Repo Selection and Code Scraping.** We first select the top 1,000 most popular pypi packages and collate the corresponding github repos where possible, similar to (Jimenez et al., 2023). These repos are then pulled down locally and filtered based on a static ast analysis determining which repositories contain type annotated code.

**2. Function Selection and Dependency Collation.** We perform a static ast analysis to filter to functions with LLM generatable argument and return type annotations. Further ast analysis then recursively identifies dependent elements (modules, functions, classes, variables, ...) across files,

builds a dependency graph, and inlines them into a base task. Finally, base tasks containing side effects or non-deterministic code such as environment variables, process calls, randomness or network requests are filtered out. See Appendix A.2 for detail on acceptable type annotations and filtering.

**3. Test Case Generation.** Using the argument type annotations we construct a LLM function calling schema that generates a diverse set of inputs. The base task code is then executed with each generated input and the result with runtime statistics are logged. This forms the test case (base task code + generated input), output (returned result or error from executed code) and statistics (runtime statistics + static ast analysis statistics).

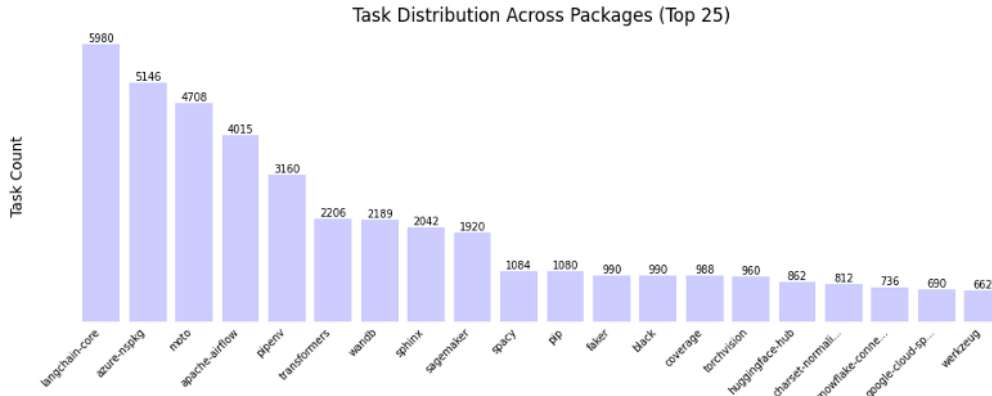


Figure 3: We observe task counts per repository to have a near logarithmic falloff. Note: manual removal of several bad offender repositories was required as they contained thousands of nearly identical functions with only url changes.

Through these stages of filtering, the original top 1,000 repositories are filtered down to the 33,950 task instances which comprise EXE. A high level breakdown of these task instances across repositories is presented in Figure 3. We note some repositories are overrepresented primarily due to being more modern (using typing) and the style of code (shorter deterministic pieces).

## 2.1 TASK FORMATION

**Model input.** The model is given a complete snippet of code alongside the input state to be executed. The model is then tasked to predict the resulting return value, or in the case that an exception is raised the model is instructed to generate an exception type and value. In practice, we prompt models with an odata json representation and use a parser to ensure valid generations. We do append one additional user reply with the parsing error if the model’s response fails to parse. Examples of input instances can be found in Appendix A.1.

**Evaluation metrics.** To evaluate a proposed solution, we use the `pass@k` metric (Chen et al., 2021), comparing the ground truth and the generated prediction as json objects (`set` and `frozenset` are sorted before conversion to json lists). If the original code produced an exception, we compare the type and message (excluding stacktrace) using a language model comparison. Examples of generated outputs can be seen in Appendix A.1.

## 2.2 FEATURES OF EXE

**Diversity of inputs and outputs.** Unlike many benchmarks focused on a particular subject matter area, a task in this eval may require a model to perform mathematical reasoning, logical inference, bit manipulation, string operations, loop execution, or to maintain multiple internal variables during computation. Furthermore, these may only form part of an algorithm that the model has to execute. Our random human inspection has uncovered algorithmic time complexities spanning from  $O(1)$  to  $O(x^n)$  and structured analysis has found tasks with code context lengths ranging from 440 to 311,000 tokens. Ensuring this broad diversity reduces the risk of hitting a local maxima and increases our opportunity to measure internal capabilities across a range of difficulties.

**Continually updatable.** Both our code collection and task input generation processes can create new tasks with minimal human oversight. Simply re-running our code collection to pull the latest commits or directing it towards an uncollected Python GitHub repository will create new task instances. Furthermore we can continue to generate new test cases for existing tasks, our test case generator automatically avoids generating seen inputs. Hence, EXE can be extended continually with new task instances, ensuring answers were not included in training corpuses of models for evaluation.

**Cost effective scalability.** With generation of new tasks requiring an average of 1,112 input tokens (batch of 15) and evaluation of tasks typically requiring 1,123 tokens, ExecEval can be generated, tested and continually updated at a fraction of the cost of human-curated benchmarks. Our initial dataset of 33,950 cases has only incurred an approximate costing of \$11 to produce and \$95 to test on.

**Long multi-step problems with smooth difficulty scaling.** We provide a continuous spectrum of task difficulties, ranging from 1-step, one-line functions to multi-file, multi-class, multi-100-step tasks. Our most complex tasks include function call depths (non-recursive) of up to 13 levels (median: 2), separate identifier counts (i.e. variable names, function names, ...) of up to 823 (median: 16) and up to 63 if statements (median: 1). This smooth scaling of difficulty allows for a more detailed measurement of model coherence along multi-step problems than what is typically seen in traditional evals. However, as language models continue to advance rapidly, even this wide range of difficulties may eventually face saturation.

To address this, we observe a mechanism inspired by the SKILL-MIX evaluation (Yu et al., 2023) that leverages the typed nature of our function selection process. This approach allows us to create even more complex tasks by chaining functions where the output type of one matches the input type of another, or by combining multiple outputs into a composite input. The number of potential new tasks can be upper bounded by  $n^2 \cdot (T_{\max})^k \cdot C$ , where  $n$  is the total number of types,  $T_{\max} = \max_{i,j} T_{i,j}$  is the maximum number of existing tasks between any two types,  $k$  is the number of functions to chain, and  $C$  is the average number of test cases per task. While this is an upper bound and the actual number of valid composite tasks would be lower due to specific type compatibility constraints, it still represents a significant expansion of our task space. We view this as an opportunity to trade some of the 'realism' of using 100% real-world code for the ability to probe the upper bounds of model capabilities. For constant compute models, this approach allows us to test their internal mechanistic capabilities in handling increasingly complex, multi-step problems. And for chain-of-thought models, it provides a test of increasingly long-term agentic coherency.

**Error prediction.** To test the full spectrum of code execution we further generate test cases designed to trigger exceptions. Many of these require in-depth analysis to see ahead of time, for example predicting an invalid array index through multiple functions. While debugging exceptions is one of the more challenging software engineering tasks, we are yet to see it commonly evaluated in benchmarks.

### 3 RESULTS

We report our evaluation results across different SOTA models alongside our findings across different task statistics below.

Table 1: EXE Pass@1 results

Model	EXE dataset (Pass@1)	Errors (Pass@1)
gpt4o	72.4	49.5
gpt4o-mini	60.9	32.0

**LLMs can execute real-world code, achieving results in-line with code generation benchmarks.** We find EXE shows similar relative model performance between models as seen in coding benchmarks such as HumanEval (Chen et al., 2021) and as seen in benchmarks requiring logical inference

such as (Lu et al., 2023). Furthermore we find a similar diversity of performance across packages as seen in agentic benchmarks such as (Jimenez et al., 2023). We show our findings in Figure 4.

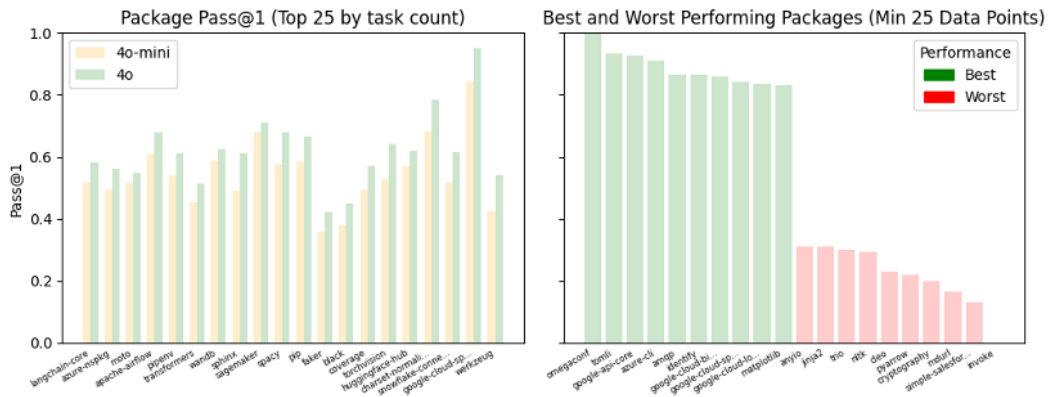


Figure 4: Left - We show the relative accuracy of different models across the top 20 packages by task count. Both the relative differences between models and the relative differences between packages are within expectations from other coding benchmarks (Jimenez et al., 2023). Right - We show the magnitude of diversity across packages (mean performance across all models).

Prior works such as Learning To Execute (Zaremba & Sutskever, 2014) and CRUX-Eval (Gu et al., 2024) have placed justifiable limitations on code complexity; removing mathematical operations, limiting line count, disallowing custom classes and only having one singular function to name a few. We hypothesised that these are no longer necessary, and to understand the true internal capabilities of a constant compute model (i.e. no Chain of Thought) we must test on real-world code, only applying limitations where forced (i.e. no arbitrary object inputs, as LLMs can't generate them). Our results as seen in table 1 provide initial evidence towards our hypothesis.

**ExecEval provides a smooth curve of task difficulties.** We set out to ensure a) our eval does not induce saturation from a bounded distribution of task difficulties, b) our eval does not induce an "AI overhang" by not having a smooth transition between difficulties and, c) the correlated factors affecting difficulty are human interpretable.

As shown in Figure 5 several task statistics such as "lines of code", "processing time" and "number of function calls" all correlate log-linearly with a model's achieved pass@1 score. These correlations provide preliminary evidence towards c) as they align with simplistic human intuition, i.e. more lines of code, more compute cycles, higher difficulty. Furthermore, we view the log-linear relationships as evidence towards b), i.e. EXE provides a smooth transition between difficulties. And finally, we view the relationships as a demonstration of difficulty being affected by factors within our control, i.e. number of function calls - providing empirical evidence towards a).

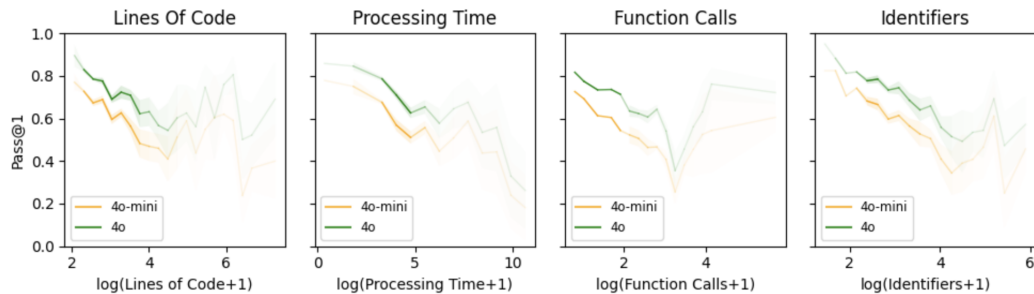


Figure 5: Pass@1 for all tasks across four of our code metrics. The shaded area represents variance, and the opacity is scaled with count of samples. Processing time is measured in microseconds.

**Stylistic coding patterns shape the metrics.** As can be seen in Figure 5 the pass@1 rate of function calls hits an elbow and then surprisingly improves as the call count increases. During our investi-

gation we found several of these occurrences, and not only with call count. These were found to be largely driven by specific coding patterns and complex tasks that LLMs excel at. We show in Figure 6 below three example tasks, and more specifically coding patterns driving this anomaly.

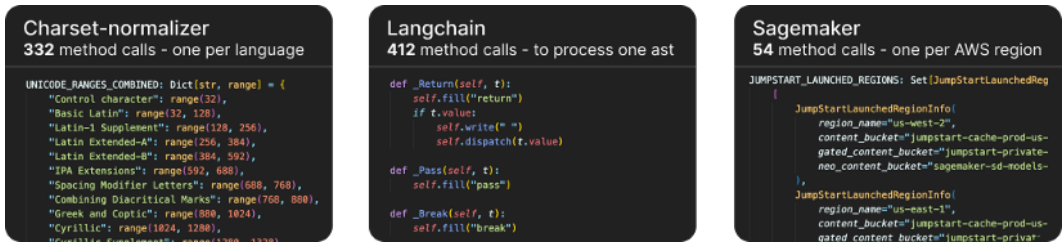


Figure 6: Three examples of high pass@1 rate tasks that contain large amounts of function calls. Left - Charset-normaliser performs 300+ function calls to define ranges of unicode characters upon initialisation; this constant has little effect on task difficulty but is used frequently and hence appears in many tasks. Middle - Langchain’s Unparser class traverses an AST and regenerates source code. The calling method in our dataset is ”add\_last\_line\_print(str) → str” which takes in code, parses it and then uses Unparse(...) to unparse it; this is a prime example of a ”directly predictable task”, i.e. one not requiring line by line code execution to predict a result. Right - Similar to Charset-normaliser, AWS’s Sagemaker has a module level constant with 10s of calls; not creating a large impact on task difficulty but frequent in its use.

**LLMs struggle with certain coding features.** As EXE contains a diverse set of tasks, we are able to observe model performance differing greatly based on coding features used in any task. To illustrate: floating point math operations such as multiplications (gpt4o: 43 mean Pass@1) significantly increase task difficulty, however bit manipulation and boolean operations only showed a minor negative impact. Iterative operations such as compound assignment operations i.e. ”i += 1” (56 Pass@1), list slicing (65 Pass@1) and list comprehensions (68 Pass@1) all increased difficulty, however for loops on (73 Pass@1) on average did not have a significant impact.

With the above metrics, and those seen in Figure 6, their mean Pass@k decreases as their count increases. To reduce the risk of our metrics being a proxy for longer problems we show the effects can still be seen below in Figure 7 after normalisation by lines of code (only lines with executable syntax tokens are counted).

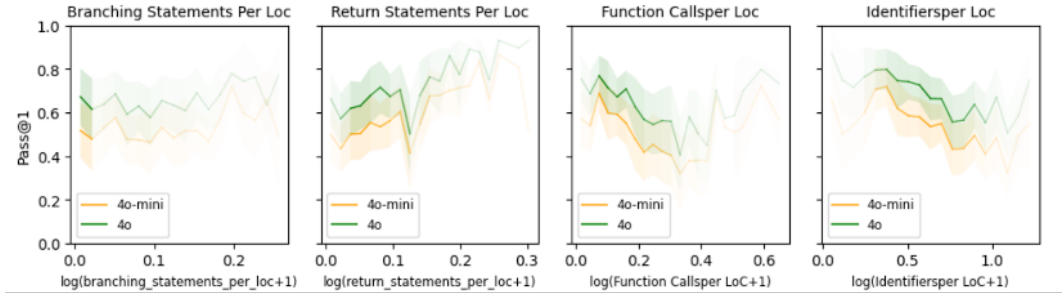


Figure 7: Pass@1for all tasks across four of our code metrics normalised by line of code count. All four of the above metrics previously showed a negative impact as they increased, interestingly we now observe branching statements having little to no impact and return statements surprisingly driving an increase in Pass@1 score. Our strong negative factors i.e. function calls and identifiers created, still are seen increasing task difficulty as they take up ever greater percentages of the task.

### 4 RELATED WORK

There is a rich history of work on evaluating language models’ abilities in reasoning, execution, and multi-step problem-solving across various domains. These efforts span from natural language processing to mathematical reasoning, and from code generation to program execution. Our work,

EXecution-Eval (EXE), builds upon this foundation while addressing key challenges in benchmark design and evaluation.

Code generation benchmarks have been the foundation of evaluating the coding abilities of language models. Works like HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) established standardized datasets for assessing code synthesis from natural language descriptions. These efforts have expanded to cover multiple programming languages (Cassano et al., 2022; Khan et al., 2023) and more complex domains such as algorithmic problem solving (Huang et al., 2023). While these benchmarks focus primarily on the task of code generation, we believe additional focus on the tasks of code execution and error prediction have been overlooked and may offer additional insight into the internal capabilities of frontier models.

The concept of "learning to execute" itself has a long history, Zaremba & Sutskever (2014) explored neural networks' ability to learn and execute simple programs. Graves et al. (2014) constructed the first Neural Turing Machines with (Kaiser & Sutskever, 2015; Reed & de Freitas, 2015; Dehghani et al., 2018) all building further into this domain. This line of research has evolved, with recent works like Bieber et al. (2020); Nye et al. (2021) and Gu et al. (2024) applying graph and language models to execute synthetic or simplistic Python programs. EXE builds upon these foundations by evaluating execution capabilities on complex, messy, real-world code from diverse GitHub repositories, providing a more challenging, scaleable and realistic test bed.

Recent trends in benchmark design have emphasized the importance of diverse, multi-step problems and agentic capabilities. Works like Jimenez et al. (2023) have introduced benchmarks that require solving real world software engineering problems while Zhou et al. (2023) has enabled evaluation of complex instruction following and performing multi-step reasoning. In the mathematical domain, benchmarks like those by Hendrycks et al. (2021) and Lu et al. (2023) have pushed models to solve intricate, multi-step problems.

The challenge of benchmark saturation and the need for continually updated evaluations has been recognized in recent works (Ott et al., 2022). Live benchmarks such as those proposed by Li et al. (2024), (Chiang et al., 2024) and Kiela et al. (2021) aim to address this issue. Skill-Mix (Yu et al., 2023) takes a novel approach, combining separate skills required to solve a problem they are able to increase task difficulty non-linearly with  $k$  skills. EXE has been inspired by both these concepts, hence the focus on enabling continual generation of new coding tasks and test cases, as well as the potential extension into chaining functions.

While many existing benchmarks use curated or synthetic datasets, EXE leverages real-world code from popular Python repositories. This approach is inspired by works like CodeNet (Puri et al., 2021) and The Stack (Kocetkov et al., 2022) which demonstrated the value of diverse, real-world data in training and evaluating language models.

## 5 EXTENSIONS

**Expanding the scope and diversity** We believe scaling EXE to include more repositories by as much as 100x would significantly reduce the noise seen in our coding metrics and provide a more resilient baseline for future frontier models. By incorporating additional Python functions — potentially using language models to predict missing type annotations — and including a diversity of other programming languages such as C++, Go and JavaScript, we believe there is even further opportunity to scale. This would offer further insights into the generalisability of a model's code understanding, pose new challenges for analysis such as pointers, macros and type-free codebases.

**Probing code execution mechanisms with simple functions** We believe there is an opportunity to align code execution with mechanistic interpretability, to gain an understanding of how constant compute language models can execute complex multi-step instructions. To illustrate, if we select the simplest function that a language model can not directly predict the outcome of, a hash function for example (one that doesn't use floating point math in this case), one requiring compute at each iteration. This would force the network to perform the computation step by step, and for a constant compute feed-forward network, layer by layer. Hence, performing a single iteration that may not lead to anything interesting, however as we increase the iteration count one by one, the model now must find a repeated circuit to perform the same computation in the later layers. For every increase it must find another circuit or a more optimal way of performing its work until it fails. We believe this



would present an interesting approach alongside standard mechanistic interpretability techniques for circuit discovery and understanding of control flow, variable tracking and computational logic at the mechanistic level.

**Breakpoint analysis for validating code execution granularly** Rather than evaluating the final return value, including multiple evaluation points within code execution may assist verification of if models are performing the step-by-step computations to reach a return value. Furthermore by inserting 'breakpoints' throughout the execution process, we can transform a single return state prediction task into numerous intermediate state prediction tasks. To illustrate, given a code snippet with a breakpoint at a specific line, a model would be tasked to determine the values of the local variables when the breakpoint is triggered. This mirrors common human debugging practices and may reveal discrepancies between final output accuracy and intermediate state understanding, offering further resistance against tasks where their final outcome can be directly predicted.

**Connection to cybersecurity threat model.** Software vulnerability research techniques are largely<sup>1</sup> enabled by the ability to predict and reason about expected program outcomes. For example, code injection, path resolution and memory buffer attacks are often found through manual human analysis; tracing inputs through the control flow, predicting output states and reasoning if there are opportunities to exploit. As EXE contains parsers such as seen in Appendix A.1 we see an opportunity to select a subset of EXE where prediction of error would imply language models have the internal capability to comprehend and aid humans with crafting vulnerabilities.

## 6 CONCLUSIONS

In this paper, we introduced EXecution-Eval (EXE), a benchmark designed to evaluate whether language models can execute real-world code. By collecting over 30,000 tasks from 1,000 popular Python repositories, EXE presents a diverse range of problems requiring computational operations such as mathematical reasoning, logical inference, and state maintenance. Our evaluations suggest that while language models demonstrate some capability in executing code, they often struggle with complex, multi-step tasks—particularly those involving many identifiers, function calls and iterative operations. Our findings indicate that although current models have limitations in accurately reasoning about and executing real-world code, they perform surprisingly well on average, prompting several opportunities extending this investigation.

EXE aims to address limitations of existing benchmarks by providing a scalable, diverse, and continually updatable framework. Its design targets a smooth difficulty scale and easy generation of new tasks with minimal human oversight with the goal to reduce the risk of models "training on the test set."

## REFERENCES

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models. August 2021.
- David Bieber, Charles Sutton, Hugo Larochelle, and Daniel Tarlow. Learning to execute programs with instruction pointer attention graph neural networks. October 2020.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. MultiPL-E: A scalable and extensible approach to benchmarking neural code generation. August 2022.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,

<sup>1</sup>Some techniques such as random fuzzing may not rely on any internal program knowledge. However, to find actionable results within realistic computational bounds, fuzzers are often augmented based on this knowledge to limit their generatable space.



- 432 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-  
433 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex  
434 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saun-  
435 ders, Christopher Hesse, Andrew N Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,  
436 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob  
437 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
438 large language models trained on code. July 2021.
- 439 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,  
440 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, and Ion Stoica.  
441 Chatbot arena: An open platform for evaluating LLMs by human preference. March 2024.
- 442
- 443 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal  
444 transformers. July 2018.
- 445
- 446 Max Fowler, David IV, Mohammed Hassan, Seth Poulsen, Matthew West, and Craig Zilles. Reeval-  
447 uating the relationship between explaining, tracing, and writing skills in cs1 in a replication study.  
448 *Computer Science Education*, 32:1–29, 06 2022. doi: 10.1080/08993408.2022.2079866.
- 449 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. October 2014.
- 450
- 451 Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I  
452 Wang. CRUXEval: A benchmark for code reasoning, understanding and execution. January 2024.
- 453
- 454 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
455 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. March  
456 2021.
- 457 Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong, Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong  
458 Shen, Chen Lin, Nan Duan, and Weizhu Chen. Competition-Level problems are effective LLM  
459 evaluators. December 2023.
- 460
- 461 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
462 Narasimhan. SWE-bench: Can language models resolve real-world GitHub issues? October  
463 2023.
- 464 Łukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. November 2015.
- 465
- 466 Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan  
467 Parvez, and Shafiq Joty. XCodeEval: A large scale multilingual multitask benchmark for code  
468 understanding, generation, translation and retrieval. March 2023.
- 469
- 470 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie  
471 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian  
472 Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina  
473 Williams. Dynabench: Rethinking benchmarking in NLP. April 2021.
- 474 Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis,  
475 Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von  
476 Werra, and Harm de Vries. The stack: 3 TB of permissively licensed source code. November  
477 2022.
- 478
- 479 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gon-  
480 zalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-Hard and  
481 BenchBuilder pipeline. June 2024.
- 482 Mike Lopez, Jacqueline Whalley, Phil Robbins, and Raymond Lister. Relationships between  
483 reading, tracing and writing skills in introductory programming. In *Proceedings of the Fourth  
484 International Workshop on Computing Education Research, ICER '08*, pp. 101–112, New  
485 York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582160. doi:  
10.1145/1404520.1404531. URL <https://doi.org/10.1145/1404520.1404531>.

- 486 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-  
487 Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of  
488 foundation models in visual contexts. October 2023.
- 489  
490 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David  
491 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Au-  
492 gustus Odena. Show your work: Scratchpads for intermediate computation with language models.  
493 November 2021.
- 494  
495 Simon Ott, Adriano Barbosa-Silva, Kathrin Blagec, Jan Brauner, and Matthias Samwald. Mapping  
496 global dynamics of benchmark creation and saturation in artificial intelligence. March 2022.
- 497  
498 Ruchir Puri, David S Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov,  
499 Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, Veronika Thost, Luca Buratti, Saurabh  
500 Pujar, Shyam Ramji, Ulrich Finkler, Susan Malaika, and Frederick Reiss. CodeNet: A large-scale  
501 AI for code dataset for learning a diversity of coding tasks. May 2021.
- 502  
503 Scott Reed and Nando de Freitas. Neural Programmer-Interpreters. November 2015.
- 504  
505 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
506 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,  
507 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W Ko-  
508 curek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda  
509 Askeel, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S Iyer, Anders An-  
510 dreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La,  
511 Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta,  
512 Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul  
513 Menezes, Arun Kirubakaran, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat,  
514 Aykut Erdem, Ayla Karakaş, B Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bo-  
515 janowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno  
516 Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron  
517 Dour, Catherine Stinson, Cedrick Argueta, and Ramírez. Beyond the imitation game: Quantifying  
518 and extrapolating the capabilities of language models. June 2022.
- 519  
520 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer  
521 Levy, and Samuel R Bowman. SuperGLUE: A stickier benchmark for general-purpose language  
522 understanding systems. May 2019.
- 523  
524 Dingli Yu, Simran Kaur, Arushi Gupta, Jonah Brown-Cohen, Anirudh Goyal, and Sanjeev Arora.  
525 Skill-Mix: a flexible and expandable family of evaluations for AI models. October 2023.
- 526  
527 Wojciech Zaremba and Ilya Sutskever. Learning to execute. October 2014.
- 528  
529 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
530 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica.  
531 Judging LLM-as-a-judge with MT-bench and chatbot arena. June 2023.
- 532  
533 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,  
534 and Le Hou. Instruction-following evaluation for large language models. November 2023.

## 531 A APPENDIX

532  
533 You may include other additional sections here.

### 534 A.1 EXAMPLE INPUT & OUTPUT

535  
536 Below is an example from the eval set. It is split into three components:

537  
538 **1. Code Task.** The function `split_email` was found to pass the type requirements, and as such  
539 all modules, classes, functions and attributes required to execute it have been recursively inlined.

540 **2. Test Case Inputs.** Based on the type definition (used for setting the function calling schema)  
 541 inputs/ output pairs have been generated with the goal of maximising diversity of control flow paths  
 542 within the function.

543 **3. Outputs.** Based on the type definition (used for setting the function calling schema) inputs/  
 544 output pairs have been generated with the goal of maximising diversity of control flow paths within  
 545 the function.  
 546

#### 547 Code

548 Note: The top 1,000 PyPI repos are used to form EXE, this function is from email-validator, rank  
 549 345  
 550

```

551 1 from typing import Optional, Tuple
552 2 import re
553 3 import unicodedata
554 4
555 5
556 6 class EmailNotValidError(ValueError):
557 7     """Parent class of all exceptions raised by this module."""
558 8     pass
559 9
560 10 class EmailSyntaxError(EmailNotValidError):
561 11     """Exception raised when an email address fails validation
562 12     because of its form."""
563 13     pass
564 14
565 15 ATEXT = r'a-zA-Z0-9_!#\$\%&\'*\+\-\/=?\^`\{\|\}\~'
566 16
567 17
568 18 def safe_character_display(c: str) -> str:
569 19     # Return safely displayable characters in quotes.
570 20     if c == '\\':
571 21         return f"\"{c}\"" # can't use repr because it escapes it
572 22     if unicodedata.category(c)[0] in ("L", "N", "P", "S"):
573 23         return repr(c)
574 24
575 25
576 26     # Construct a hex string in case the unicode name doesn't exist.
577 27     if ord(c) < 0xFFFF:
578 28         h = f"U+{ord(c):04x}".upper()
579 29     else:
580 30         h = f"U+{ord(c):08x}".upper()
581 31
582 32     # Return the character name or, if it has no name, the hex
583 33     string.
584 34     return unicodedata.name(c, h)
585 35
586 36
587 37 ATEXT_RE = re.compile('[' + ATEXT + ']') # ATEXT plus dots
588 38
589 39
590 40 def check_unsafe_chars(s: str, allow_space: bool = False) -> None:
591 41     # Check for unsafe characters or characters that would make the
592 42     string
593 43     # invalid or non-sensible Unicode.
594 44     bad_chars = set()
595 45     for i, c in enumerate(s):
596 46         category = unicodedata.category(c)
597 47         if category[0] in ("L", "N", "P", "S"):
598 48             # Letters, numbers, punctuation, and symbols are permitted

```

```

594 49     pass
595 50     elif category[0] == "M":
596 51         # Combining character in first position would combine with
597 52             something
598 53         # outside of the email address if concatenated, so they
599 54             are not safe.
600 55         # We also check if this occurs after the @-sign, which
601 56             would not be
602 57         # sensible because it would modify the @-sign.
603 58         if i == 0:
604 59             bad_chars.add(c)
605 60     elif category == "Zs":
606 61         # Spaces outside of the ASCII range are not specifically
607 62             disallowed in
608 63         # internationalized addresses as far as I can tell, but
609 64             they violate
610 65         # the spirit of the non-internationalized specification
611 66             that email
612 67         # addresses do not contain ASCII spaces when not quoted.
613 68             Excluding
614 69         # ASCII spaces when not quoted is handled directly by the
615 70             atom regex.
616 71         #
617 72         # In quoted-string local parts, spaces are explicitly
618 73             permitted, and
619 74         # the ASCII space has category Zs, so we must allow it
620 75             here, and we'll
621 76         # allow all Unicode spaces to be consistent.
622 77         if not allow_space:
623 78             bad_chars.add(c)
624 79     elif category[0] == "Z":
625 80         # The two line and paragraph separator characters (in
626 81             categories Zl and Zp)
627 82         # are not specifically disallowed in internationalized
628 83             addresses
629 84         # as far as I can tell, but they violate the spirit of the
630 85             non-internationalized
631 86         # specification that email addresses do not contain line
632 87             breaks when not quoted.
633 88         bad_chars.add(c)
634 89     elif category[0] == "C":
635 90         # Control, format, surrogate, private use, and unassigned
636 91             code points (C)
637 92         # are all unsafe in various ways. Control and format
638 93             characters can affect
639 94         # text rendering if the email address is concatenated with
640 95             other text.
641 96         # Bidirectional format characters are unsafe, even if used
642 97             properly, because
643 98         # they cause an email address to render as a different
644 99             email address.
645 100        # Private use characters do not make sense for publicly
646 101            deliverable
647 102        # email addresses.
648 103        bad_chars.add(c)
649 104    else:
650 105        # All categories should be handled above, but in case
651 106            there is something new
652 107        # to the Unicode specification in the future, reject all
653 108            other categories.
654 109        bad_chars.add(c)
655 110    if bad_chars:
656 111        raise EmailSyntaxError("The email address contains unsafe
657 112            characters: ")

```

```

648 90         + ", ".join(safe_character_display(c)
649         for c in sorted(bad_chars)) + ".")
650 91
651 92
652 93     def split_email(email: str) -> Tuple[Optional[str], str, str, bool]:
653 94         # Return the display name, unescaped local part, and domain part
654 95         # of the address, and whether the local part was quoted. If no
655 96         # display name was present and angle brackets do not surround
656 97         # the address, display name will be None; otherwise, it will be
657 98         # set to the display name or the empty string if there were
658 99         # angle brackets but no display name.
659 100
660 101
661 102         # Typical email addresses have a single @-sign and no quote
662 103         # characters, but the awkward "quoted string" local part form
663 104         # (RFC 5321 4.1.2) allows @-signs and escaped quotes to appear
664 105         # in the local part if the local part is quoted.
665 106
666 107
667 108         # A 'display name <addr>' format is also present in MIME
668 109         # messages
669 110         # (RFC 5322 3.4) and this format is also often recognized in
670 111         # mail UIs. It's not allowed in SMTP commands or in typical web
671 112         # login forms, but parsing it has been requested, so it's done
672 113         # here as a convenience. It's implemented in the spirit but not
673 114         # the letter of RFC 5322 3.4 because MIME messages allow
674 115         # newlines
675 116         # and comments as a part of the CFWS rule, but this is typically
676 117         # not
677 118         # allowed in mail UIs (although comment syntax was requested
678 119         # once too).
679 120         #
680 121         # Display names are either basic characters (the same basic
681 122         # characters
682 123         # permitted in email addresses, but periods are not allowed and
683 124         # spaces
684 125         # are allowed; see RFC 5322 Appendix A.1.2), or or a quoted
685 126         # string with
686 127         # the same rules as a quoted local part. (Multiple quoted
687 128         # strings might
688 129         # be allowed? Unclear.) Optional space (RFC 5322 3.4 CFWS) and
689 130         # then the
690 131         # email address follows in angle brackets.
691 132         #
692 133         # We assume the input string is already stripped of leading and
693 134         # trailing CFWS.
694 135
695 136
696 137     def split_string_at_unquoted_special(text: str, specials: Tuple[
697 138         str, ...]) -> Tuple[str, str]:
698 139         # Split the string at the first character in specials (an @-
699 140         # sign
700 141         # or left angle bracket) that does not occur within quotes
701 142         # and
702 143         # is not followed by a Unicode combining character.
703 144         # If no special character is found, raise an error.
704 145         inside_quote, escaped, left_part = False, False, ""
705 146         for i, c in enumerate(text):
706 147             # < plus U+0338 (Combining Long Solidus Overlay)
707 148             # normalizes to
708 149             # U+226E (Not Less-Than), and it would be confusing to
709 150             # treat
710 151             # the < as the start of "<email>" syntax in that case.
711 152             Likewise,

```

```

702 137     # if anything combines with an @ or ", we should probably
703 138     not
704 138     # treat it as a special character.
705 139     if unicodedata.normalize("NFC", text[i:])[0] != c:
706 140         left_part += c
707 141
708 142
709 143     elif inside_quote:
710 144         left_part += c
711 145         if c == '\\\' and not escaped:
712 146             escaped = True
713 147         elif c == '\"' and not escaped:
714 148             # The only way to exit the quote is an unescaped
715 149             quote.
716 150             inside_quote = False
717 151             escaped = False
718 152         else:
719 153             escaped = False
720 154     elif c == '\"':
721 155         left_part += c
722 156         inside_quote = True
723 157     elif c in specials:
724 158         # When unquoted, stop before a special character.
725 159         break
726 160     else:
727 161         left_part += c
728 162
729 163
730 164     if len(left_part) == len(text):
731 165         raise EmailSyntaxError("An email address must have an @-
732 166         sign.")
733 167
734 168     right_part = text[len(left_part):] # The right part is
735 169     whatever is left.
736 170
737 171     return left_part, right_part
738 172
739 173     def unquote_quoted_string(text: str) -> Tuple[str, bool]:
740 174         # Remove surrounding quotes and unescape escaped backslashes
741 175         # and quotes. Escapes are parsed liberally. I think only
742 176         backslashes
743 177         # and quotes can be escaped but we'll allow anything to be.
744 178         quoted, escaped, value = False, False, ""
745 179         for i, c in enumerate(text):
746 180             if quoted:
747 181                 if escaped:
748 182                     value += c
749 183                     escaped = False
750 184                 elif c == '\\\' :
751 185                     escaped = True
752 186                 elif c == '\"':
753 187                     if i != len(text) - 1:
754 188                         raise EmailSyntaxError("Extra character(s)
755 189                         found after close quote: "
760 190                         + ", ".join(
761 191                             safe_character_display
762 192                             (c) for c in text[i +
763 193                             1:]))
764 194                     break
765 195                 else:
766 196                     value += c
767 197             elif i == 0 and c == '\"':

```

```

756 193         quoted = True
757 194     else:
758 195         value += c
759 196
760 197
761 198     return value, quoted
762 199
763 200
764 201     # Split the string at the first unquoted @-sign or left angle
765 202     left_part, right_part = split_string_at_unquoted_special(email,
766 203         ("@", "<"))
767 204
768 205     # If the right part starts with an angle bracket, then the left
769 206     part
770 207     # is a display name and the rest of the right part up to the
771 208     # final right angle bracket is the email address, .
772 209     if right_part.startswith("<"):
773 210         # Remove space between the display name and angle bracket.
774 211         left_part = left_part.rstrip()
775 212
776 213     # Unquote and unescape the display name.
777 214     display_name, display_name_quoted = unquote_quoted_string(
778 215         left_part)
779 216
780 217     # Check that only basic characters are present in a non-
781 218     quoted display name.
782 219     if not display_name_quoted:
783 220         bad_chars = {
784 221             safe_character_display(c)
785 222             for c in display_name
786 223             if (not ATEXT_RE.match(c) and c != ' ') or c == '.'
787 224         }
788 225     if bad_chars:
789 226         raise EmailSyntaxError("The display name contains
790 227         invalid characters when not quoted: " + ", ".
791 228         join(sorted(bad_chars)) + ".")
792 229
793 230     check_unsafe_chars(display_name, allow_space=True) # Check
794 231     for other unsafe characters.
795 232
796 233     # Check that the right part ends with an angle bracket
797 234     # but allow spaces after it, I guess.
798 235     if ">" not in right_part:
799 236         raise EmailSyntaxError("An open angle bracket at the start
800 237         of the email address has to be followed by a close
801 238         angle bracket at the end.")
802 239     right_part = right_part.rstrip(" ")
803 240     if right_part[-1] != ">":
804 241         raise EmailSyntaxError("There can't be anything after the
805 242         email address.")
806 243
807 244     # Remove the initial and trailing angle brackets.
808 245     addr_spec = right_part[1:].rstrip(">")
809
810     # Split the email address at the first unquoted @-sign.
811     local_part, domain_part = split_string_at_unquoted_special(
812         addr_spec, ("@",))

```



```

810 246
811 247
812 248     # Otherwise there is no display name. The left part is the local
813 249     # part and the right part is the domain.
814 250     else:
815 251         display_name = None
816 252         local_part, domain_part = left_part, right_part
817 253
818 254
819 255     if domain_part.startswith("@"):
820 256         domain_part = domain_part[1:]
821 257
822 258
823 259     # Unquote the local part if it is quoted.
824 260     local_part, is_quoted_local_part = unquote_quoted_string(
825 261         local_part)
826 262
827 263     return display_name, local_part, domain_part,
828         is_quoted_local_part

```

### Test Case Inputs

```

828
829
830 1  [
831 2  {
832 3  "input": [{"simple@example.com"}, {}],
833 4  "output": [null,"simple","example.com", false],
834 5  },
835 6  {
836 7  "input": [{"user+name@sub.domain.com"}, {}],
837 8  "output": [null,"user+name","sub.domain.com", false],
838 9  },
839 10 {
840 11 "input": [{"user.name@domain.co.uk"}, {}],
841 12 "output": [null,"user.name","domain.co.uk", false],
842 13 },
843 14 {
844 15 "input": [{"\"quoted@local\"@example.com"}, {}],
845 16 "output": [null,"quoted@local","example.com", true],
846 17 },
847 18 {
848 19 "input": [{"display name <user@domain.com>"}, {}],
849 20 "output": ["display name","user","domain.com", false],
850 21 },
851 22 {
852 23 "input": [{"user@localhost"}, {}],
853 24 "output": [null,"user","localhost", false],
854 25 },
855 26 {
856 27 "input": [{"user@[IPv6:2001:db8::1]"}, {}],
857 28 "output": [null,"user","[IPv6:2001:db8::1]", false],
858 29 },
859 30 {
860 31 "input": [{"\"escaped\\\"quote\"@example.com"}, {}],
861 32 "output": [null,"escaped\"quote","example.com", true],
862 33 },
863 34 {
864 35 "input": [{"user.name@longsubdomain.example.com"}, {}],
865 36 "output": [null,"user.name","longsubdomain.example.com", false],
866 37 },
867 38 {
868 39 "input": [{"very.common@example.com"}, {}],
869 40 "output": [null,"very.common","example.com", false],
870 41 },
871 42 {

```

```

864 43 "input": [{"user@domain-with-dash.com"}, {}],
865 44 "output": [null,"user","domain-with-dash.com", false],
866 45 },
867 46 {
868 47 "input": [{"user@123.123.123.123"}, {}],
869 48 "output": [null,"user","123.123.123.123", false],
870 49 },
871 50 {
872 51 "input": [{"\"much.more unusual\"@example.com"}, {}],
873 52 "output": [null,"much.more unusual","example.com", true],
874 53 },
875 54 {
876 55 "input": [{"user@xn--exmple-cua.com"}, {}],
877 56 "output": [null,"user","xn--exmple-cua.com", false],
878 57 },
879 58 {
880 59 "input": [{"user@domain_with_underscore.com"}, {}],
881 60 "output": [null,"user","domain_with_underscore.com", false],
882 61 }
883 62 ]

```

### Generated Outputs

```

884 1 [
885 2 {
886 3 "input": [{"simple@example.com"}, {}],
887 4 "output": [null,"simple","example.com", false],
888 5 "prediction": [null,"simple","example.com",false],
889 6 "result": true,
890 7 "answer_tokens": {"completion": 18,"prompt": 4610,"total": 4628}
891 8 },
892 9 {
893 10 "input": [{"user+name@sub.domain.com"}, {}],
894 11 "output": [null,"user+name","sub.domain.com", false],
895 12 "prediction": [null,"user+name","sub.domain.com",false],
896 13 "result": true,
897 14 "answer_tokens": {"completion": 21,"prompt": 4614,"total": 4635}
898 15 },
899 16 {
900 17 "input": [{"user.name@domain.co.uk"}, {}],
901 18 "output": [null,"user.name","domain.co.uk", false],
902 19 "prediction": [null,"user.name","domain.co.uk",false],
903 20 "result": true,
904 21 "answer_tokens": {"completion": 20,"prompt": 4613,"total": 4633}
905 22 },
906 23 {
907 24 "input": [{"\"quoted@local\"@example.com"}, {}],
908 25 "output": [null,"quoted@local","example.com", true],
909 26 "prediction": [null,"quoted@local","example.com",true],
910 27 "result": false,
911 28 "answer_tokens": {"completion": 20,"prompt": 4615,"total": 4635}
912 29 },
913 30 {
914 31 "input": [{"display name <user@domain.com>"}, {}],
915 32 "output": [null,"display name","user","domain.com", false],
916 33 "prediction": [null,"display name","user","domain.com",false],
917 34 "result": true,
918 35 "answer_tokens": {"completion": 19,"prompt": 4615,"total": 4634}
919 36 },
920 37 {
921 38 "input": [{"user@localhost"}, {}],
922 39 "output": [null,"user","localhost", false],
923 40 "prediction": [null,"user","localhost",false],
924 41 "result": true,
925 42 "answer_tokens": {"completion": 17,"prompt": 4610,"total": 4627}

```

```

918 43 },
919 44 {
920 45   "input": [{"user@[IPv6:2001:db8::1]"}, {}],
921 46   "output": [null, "user", "[IPv6:2001:db8::1]", false],
922 47   "prediction": "EmailSyntaxError: An email address must have an @-
923 48     sign.",
924 49   "result": false,
925 50   "answer_tokens": {"completion": 24, "prompt": 4620, "total": 4644}
926 51 },
927 52 {
928 53   "input": [{"\"escaped\\\"quote\\\"@example.com"}, {}],
929 54   "output": [null, "escaped\"quote", "example.com", true],
930 55   "prediction": [null, "escaped\"quote", "example.com", true],
931 56   "result": false,
932 57   "answer_tokens": {"completion": 20, "prompt": 4615, "total": 4635}
933 58 },
934 59 {
935 60   "input": [{"user.name@longsubdomain.example.com"}, {}],
936 61   "output": [null, "user.name", "longsubdomain.example.com", false],
937 62   "prediction": [null, "user.name", "longsubdomain.example.com", false],
938 63   "result": true,
939 64   "answer_tokens": {"completion": 22, "prompt": 4615, "total": 4637}
940 65 },
941 66 {
942 67   "input": [{"very.common@example.com"}, {}],
943 68   "output": [null, "very.common", "example.com", false],
944 69   "prediction": [null, "very.common", "example.com", false],
945 70   "result": true,
946 71   "answer_tokens": {"completion": 19, "prompt": 4611, "total": 4630}
947 72 },
948 73 {
949 74   "input": [{"user@domain-with-dash.com"}, {}],
950 75   "output": [null, "user", "domain-with-dash.com", false],
951 76   "prediction": [null, "user", "domain-with-dash.com", false],
952 77   "result": true,
953 78   "answer_tokens": {"completion": 21, "prompt": 4614, "total": 4635}
954 79 },
955 80 {
956 81   "input": [{"user@123.123.123.123"}, {}],
957 82   "output": [null, "user", "123.123.123.123", false],
958 83   "prediction": [null, "user", "123.123.123.123", false],
959 84   "result": true,
960 85   "answer_tokens": {"completion": 23, "prompt": 4616, "total": 4639}
961 86 },
962 87 {
963 88   "input": [{"\"much.more unusual\\\"@example.com"}, {}],
964 89   "output": [null, "much.more unusual", "example.com", true],
965 90   "prediction": [null, "much.more unusual", "example.com", true],
966 91   "result": true,
967 92   "answer_tokens": {"completion": 20, "prompt": 4615, "total": 4635}
968 93 },
969 94 {
970 95   "input": [{"user@xn--exmple-cua.com"}, {}],
971 96   "output": [null, "user", "xn--exmple-cua.com", false],
972 97   "prediction": [null, "user", "xn--exmple-cua.com", false],
973 98   "result": true,
974 99   "answer_tokens": {"completion": 24, "prompt": 4617, "total": 4641}
975 100 },
976 101 {
977 102   "input": [{"user@domain_with_underscore.com"}, {}],
978 103   "output": [null, "user", "domain_with_underscore.com", false],
979 104   "prediction": "EmailSyntaxError: The email address contains unsafe
980 105     characters: 'U+005F'.",
981 106   "result": false,

```

```
972 106     "answer_tokens": {"completion": 28, "prompt": 4614, "total": 4642}
973 107     }
974 108 ]
```

## 977 A.2 ACCEPTABLE TYPES & FILTERING CRITERIA

978 **Acceptable types.** To find functions where the inputs and outputs are LLM generat-  
979 able, we recursively parse both arguments and return types as ast objects i.e. for  
980 list[tuple[str, False]] we first check list is an acceptable type, then recurse down  
981 into tuple, following that we then check str and finally we check False. False isn't an ac-  
982 ceptable type but it is an acceptable constant and hence accepted. Note: certain acceptable types and  
983 constants are not allowed as return values, i.e. None is not an accepted return constant

```
984 acceptable_types = { 'int', 'str', 'float', 'bool', 'none', 'list', 'dict',  
985 'tuple', 'set', 'datetime.date', 'date', 'literal', 'optional', 'union',  
986 'sequence', 'iterable', 'frozenset', 'mapping' }
```

```
987 acceptable_constants = { 'ellipsis', True, False, None }
```

989 **Filtering functions.** When filtering functions we maintain four separate block lists, 1) a list of  
990 banned imports (including direct and aliases), 2) a list of banned functions (some common libraries  
991 have a limited set of non-deterministic methods, we don't want to fully exclude them), 3) a list of  
992 banned variables (some variables such as `__version__` are likely to be environment based), 4) a  
993 list of banned repos (some repos from cloud providers provide thousands of near identical methods  
994 with different urls, we remove these as they are not a valuable contribution to the evaluation).

995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025