

Program Semantic Inequivalence Game with Large Language Models

Anonymous ACL submission

Abstract

Large Language Models (LLMs) can achieve strong performance on everyday coding tasks, but they can fail on complex tasks that require non-trivial reasoning about program semantics. Finding training examples to teach LLMs to solve these tasks can be challenging.

In this work, we explore a method to synthetically generate code reasoning training data based on a **semantic inequivalence game (SInQ)**: a **generator** agent creates program variants that are semantically distinct, derived from a dataset of real-world programming tasks, while an **evaluator** agent has to identify input examples for which they behave differently. The agents train each other semi-adversarially, improving their ability to understand the underlying logic of code.

We evaluated our approach on multiple code generation and understanding benchmarks, including cross-language **vulnerability detection** (Lu et al., 2021), where our method improves vulnerability detection in C/C++ code despite being trained exclusively on Python code, and the challenging **Python builtin identifier swap** benchmark (Miceli Barone et al., 2023), showing that whereas modern LLMs still struggle with this benchmark, our approach yields substantial improvements.

We release the code needed to replicate the experiments, as well as the generated synthetic data, which can be used to fine-tune LLMs.

1 Introduction

Assistants based on Large Language Models (LLMs) are widely used by programmers for coding tasks. While they perform well on common tasks, they still struggle with non-trivial reasoning about program semantics (Miceli Barone et al., 2023; Maveli et al., 2025). This limitation can lead to subtle bugs or prevent the detection of preexisting vulnerabilities and adversarial backdoors (Dinh

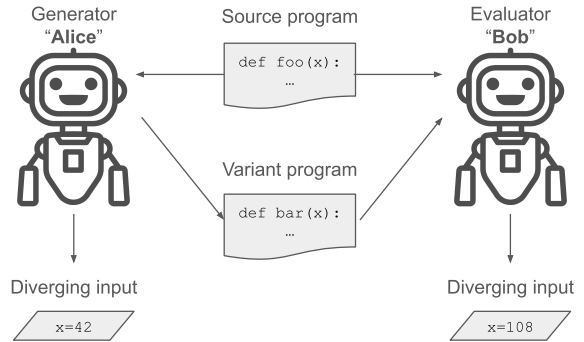


Figure 1: Semantic inequivalence game: Alice receives a source program P and generates a variant program Q and a diverging input. Bob receives P and Q and generates another diverging input.

et al., 2023; Dou et al., 2024), ultimately compromising the safety and security of generated code (Wang et al., 2024; Mohsin et al., 2024).

LLMs' code generation and understanding capabilities are typically improved by fine-tuning on a mixture of human-annotated and synthetically generated data. However, human annotation is expensive and fails to cover many non-trivial scenarios. Typical synthetic generation approaches rely on LLMs to generate coding problem statements along with corresponding solutions and unit tests, then validate solutions by executing them against the tests. While this allows for large-scale dataset creation, it may still provide limited coverage of problem types and introduce noise, as unit tests do not always align well with problem statements, particularly in edge cases.

Self-play involves training AI agents by pitting them against each other in adversarial games, incentivizing them to discover and defend against unusual scenarios. This approach has enabled AI systems to achieve human-level or even superhuman performance in games such as Go (Silver et al., 2016, 2017b), Chess (Silver et al., 2017a), Dota 2 (OpenAI et al., 2019), StarCraft II (Arulkumaran

et al., 2019), and even social games involving dialogue, such as Diplomacy (FAIR, 2022). However, self-play methods typically require external engines to enforce game rules and compute scores, making them challenging to apply to open-ended tasks like coding. Recreational competitive coding environments such as CROBOTS¹ are overly domain-specific and impose strict performance limits, making them unsuitable for training agents in general code reasoning. We are aware of only one work, Zhao et al. (2025), concurrent to our own, which uses self-play to train LLMs for arbitrary code generation, while Dong and Ma (2025) use a similar approach for theorem proving.

In this work, we introduce a game based on program semantic inequivalence designed to train agents in code reasoning across arbitrary domains. By design, this game has no theoretical performance cap. We use it to train LLMs through self-play, demonstrating significant performance improvements on challenging tasks.

2 Proposed method

Our approach involves two LLM agents engaging in a game where the **generator** agent, "Alice," creates challenging code understanding problems for the **evaluator** agent, "Bob," to solve. Alice's goal is to deceive Bob into making mistakes, requiring her to generate difficult instances. However, Alice must also provide solutions, meaning the instances cannot be unsolvable or excessively difficult. We train Alice to become more effective at misleading Bob and Bob to become better at resisting deception, encouraging both agents to develop a deeper understanding of program semantics.

Our approach is based on the **semantic equivalence** of programs, or more specifically, **semantic inequivalence**. This allows for precise verification of solutions, unlike problems based on natural language specifications or unit tests, which offer only partial coverage. Moreover, it is fundamentally linked to computability theory through reductions to Rice's theorem and the Halting problem.

In practical applications, reasoning about program equivalence and inequivalence is valuable for identifying bugs and security vulnerabilities introduced during code refactoring.

¹<https://github.com/tpoindex/crobots>

2.1 The semantic inequivalence game

Two programs P and Q are semantically equivalent if, for any given input x , either they both halt producing the same output y or neither halts. Determining semantic equivalence is a fundamental problem in program verification and compiler design, but automatic proving equivalence between arbitrary programs is complicated since popular programming languages, such as Java and Python, are defined through natural language specifications or reference implementations rather than formal semantics. Even when formal semantics exist for a subset of a language, automatically generating fully machine-checkable equivalence proofs for non-trivial code is quite challenging even for expert human programmers. We sidestep this issue by defining a program understanding game that focuses solely on inequivalent programs.

We define the **Semantic Inequivalence Game** as the following one-shot interaction between two players: the **generator**, "Alice," and the **evaluator**, "Bob":

1. Alice receives a program P and generates another program Q , which has to be inequivalent to P , along with a diverging input x such that $P(x) \neq Q(x)$.
2. The diverging input is verified by executing both programs on it. If $P(x) = Q(x)$, Alice loses.
3. Bob receives P and Q and attempts to produce a diverging input \hat{x} (which may or may not be the same as x). If Bob correctly identifies a diverging input, he wins and Alice loses; otherwise, Bob loses and Alice wins.

Alice's objective is to generate challenging instances for Bob, while Bob's goal is to solve them. In this game, correctness can be verified simply by executing the programs on the provided diverging inputs, eliminating the need for generating and verifying formal proofs.

Both agents are trained iteratively through repeated play. The source programs P provided to Alice are sampled from a dataset, such as a collection of short, self-contained programming exercises spanning a variety of tasks (e.g., MBPP (Austin et al., 2021)). This ensures that the game stays **grounded** to real-world coding problems. If Alice were allowed to generate both P and Q , she

might develop an incentive to produce unusual, obfuscated code that might not contribute to Bob’s general reasoning skills.

To approximate non-termination detection, we impose a randomized time limit that significantly exceeds the typical runtime of source programs. This prevents Alice from exploiting a fixed time limit, for example, by generating a program Q that loops for a predetermined duration before returning the same output as P .²

Example 1.³ 1. Alice receives program P :

```
def fib(n):
    if n <= 0:
        return 0
    elif n == 1:
        return 1
    return fib(n - 1) + fib(n - 2)
```

and returns program Q :

```
def fib(n):
    if n == 0:
        return 0
    elif n == 1:
        return 1
    return fib(n - 1) + fib(n - 2)
```

together with diverging input x :

```
{"n": -1}
```

2. Both programs are run in a sandbox with input x , which results in P returning 0 while Q recurs until it triggers either the recursion limit of the Python interpreter or the randomized time limit, proving that x is indeed a diverging input.

3. Bob receives both P and Q a generates a possibly different diverging input \hat{x} , e.g.:

```
{"n": -2}
```

P and Q are evaluated on input \hat{x} , proving that this is also a diverging input, therefore, Bob wins this round.

Unlike games such as Go or Chess, where perfect play is theoretically possible, the semantic inequivalence game has no strict performance cap: given an infinite time limit, Bob’s task is undecidable (see Appendix A). This implies that in principle the agents can learn arbitrarily complex coding logic while remaining grounded in a dataset of practical, real-world programming problems.

The semantic inequivalence game is entirely adversarial and essentially a zero-sum game, provided

²The time limit is randomized to discourage Alice from gaming the system by crafting artificial delays, which could lead to uninteresting cases.

³This example is artificial, please refer to section 3.3 for an analysis of a real example from the MBPP dataset.

that Alice generates only valid outputs. In some cases, modifying the game to be positive-sum may be beneficial, both to facilitate integration with supervised fine-tuning (SFT) and to prevent degenerate strategies (e.g., Alice generating excessively difficult cryptographic puzzles). We discuss these considerations further in Appendix B.

2.2 Implementation with Supervised Fine-tuning and Difficulty Targeting

The semantic inequivalence game, as defined above, is well-suited for reinforcement learning, however, reinforcement learning was not available on the OpenAI API at the time of our experiments, therefore, we devised the following rejection sampling fine-tuning implementation, with explicit difficulty supervision.

When we present the program pair (P, Q) to Bob, we can sample N evaluation responses and define the difficulty of the pair based on the number of correct assessments:

$$d(P, Q, Bob) = 10 \left(1 - \frac{N_{\text{correct}}}{N} \right)$$

For example, if Bob can solve the pair (P, Q) 40% of the time, the difficulty of this instance is 6.

During generation, we ask Alice to generate a program with a specific target difficulty d_t , usually set to the maximum value of 10 (though in some cases, we may set it to a lower value, making the game positive-sum; see Appendix B).

Let:

$$I = \text{Template}_{\text{Alice}}(P, d_t)$$

$$O = \text{Alice}(I)$$

$$(CoT, Q, x) = \text{Extractor}_{\text{Alice}}(O)$$

If Alice’s output is invalid, we discard it. Otherwise, we evaluate it with Bob to estimate its actual difficulty. We then create a new SFT training example for Alice TE_{Alice} by substituting the estimated difficulty (rounded to the nearest integer) with the target difficulty in the input. That is, we treat the actual generated program’s difficulty as if it were the target difficulty:

$$TE_{\text{Alice}} := (\text{Template}_{\text{Alice}}(P, d(P, Q, Bob)), O)$$

We can generate one or more training examples for Alice from the P programs in the source program set, then batch-train Alice, for instance, using OpenAI’s fine-tuning API with the chat LLM format. Here, the input is encoded as the "user" message

and the output as the "assistant" message, with the same "system" message used during generation. The loss is minimized only on the "assistant" message.

We can continue to extend this dataset across multiple generations of Alice, as long as Bob remains unchanged. Once we update Bob (using rejection sampling SFT on its own successful input-output pairs), we need to recompute all the difficulty estimates for the programs in Alice’s dataset, as Bob is presumably stronger. We have found it beneficial to train Alice for many epochs before training a new Bob. Initially, Alice tends to generate examples that are too easy for Bob (especially when both Alice and Bob are based on the same LLM). Ideally, we would continue training Alice until convergence before each new Bob training run.

Since Alice’s initial examples are often very easy for Bob (with difficulty zero for most), using all of them as training examples would overwhelm Alice’s training set with unhelpful instances. This could be detrimental, as we would minimize the loss on tokens of programs that we don’t want Alice to generate. To address this, we bias the training set by selecting all hard examples, defined as $d(P, Q, Bob) \geq 5$, i.e., the examples that fool Bob at least 50% of the time, plus a fraction of the remaining examples (20% of the number of hard examples), sampled without replacement by going through difficulty bins in a round-robin fashion.

We also explicitly train Alice to predict instance difficulty by including training examples in the following format:

$$TE_{diff} := (Template_{Alice}(P, "Any"), \\ O, \\ Template_{diff_{in}}, \\ Template_{diff_{out}}(d(P, Q, Bob)))$$

where the target difficulty in the first "user" message is replaced by the string "Any", and the first "assistant" message contains Alice’s self-generated output instance. The second "user" message provides an instruction to predict the difficulty of the instance, and the second "assistant" message contains the actual difficulty. We minimize the loss only on the second "assistant" message. This part of the dataset is also biased towards hard examples, but we set the number of easy examples at 50%, as we are not minimizing the loss on the tokens of easy instances but only on their difficulty

prediction. Therefore, including these examples is unlikely to be detrimental.

Refer to Appendix C for all the templates used to interact with the LLM.

3 Experiments

3.1 Training

We run our main set of experiments using OpenAI gpt-4o-mini-2024-07-18 as the base LLM for both Alice and Bob. In order to train our agents, we use the training portion of MBPP as our source set of programs, using only the code field from each source example. We perform additional experiments using OpenAI gpt-4.1-nano-2025-04-14 as the base LLM.

Both Alice and Bob are instructed to produce output in markdown format, using markdown sections to distinguish their CoT traces from the final outputs, which are embedded in Python code blocks. We sample from the models with a temperature of 1.0 and top_p = 0.7, generating $N = 10$ samples per query. We use the Mistune markdown parser⁴ followed by the Python ast parser/unparser. This step extracts, syntactically validates, and normalizes the outputs⁵. We then semantically check the diverging inputs against the pairs of source and generated programs by executing them within a test harness, using a randomized time limit, uniformly sampled between 2.5 and 5.5 seconds, to discourage Alice from generating instances that rely on race conditions.

We train the models via SFT with difficulty targeting (always set to 10) and difficulty prediction, as described in Section 2.2, using the default hyperparameters of the OpenAI fine-tuning platform.

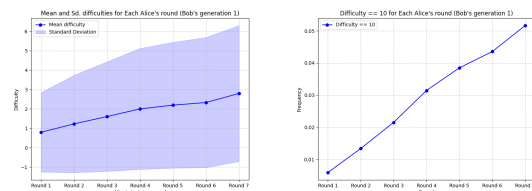


Figure 2: Instance difficulty with respect to an untrained Bob, plotted against the number of Alice’s training rounds, gpt-4-mini models. Left: mean and standard deviation. Right: fraction of instances with maximum difficulty.

⁴<https://mistune.lepture.com/en/latest/>

⁵Parsing and then unparsing the Python code with ast removes comments or unusual indentation styles, preventing trivial non-semantic attacks.

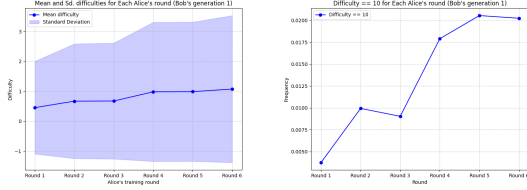


Figure 3: Instance difficulty with respect to an untrained Bob, plotted against the number of Alice’s training rounds, gpt-4.1-nano models. Left: mean and standard deviation. Right: fraction of instances with maximum difficulty.

For initial set of the experiments on the gpt-4-mini models, we perform 7 batched training rounds for Alice, followed by a single training round for Bob. This was due to time and financial constraints; ideally, we would perform training rounds for Alice until convergence of the average instance difficulty before performing one training round for Bob. For the additional set of experiments on the gpt-4.1-nano models, we perform 6 batched training rounds for Alice, which is sufficient for convergence, followed by a single training round for Bob. We report the difficulty curves in Figure 2 and 3.

Each of Alice’s training runs starts from the baseline LLM checkpoint rather than the previous fine-tuned checkpoint, but we accumulate instances to be used as training examples between rounds. Since Bob does not change between Alice’s training rounds, the difficulty estimate of each instance remains valid. However, if we were to perform additional training rounds for Alice after training Bob, we would have to either discard the training set or re-estimate the difficulty of each instance by evaluating it with the new Bob.

We consider the fine-tuned Bob to be our final model, which we use for evaluation.

We also attempted a set of experiments based on the Alibaba Qwen3 Thinking LLM, specifically the smallest non-quantized version available on Hugging Face Qwen3-4B-Thinking-2507. We performed generation using the recommended hyperparameters, with full context length (32768 tokens). We used a rank-stabilized LoRA adapter with rank = 128, alpha = 16, dropout probability = 0.1, epochs = 5, maximum learning rate = $2e - 4$. We did not use difficulty prediction examples, since according to the Qwen3 fine-tuning guidelines, CoT traces should not be included in assistant messages in multi-turn conversations ex-

Source programs	Untrained	Trained
MBPP Train	75.99	86.98
MBPP Test	88.37	91.67

Table 1: Percentage of semantic-inequivalence game instances solved by Bob, without or with training, using gpt-4-mini.

cept on the last one. This sets of experiments lead to unstable results, with the proportion of examples with maximum difficulty increasing with each Alice’s round, as expected, but the mean difficulty decreasing. Training curves are reported in Figure 4. This may be due to the small size of the model or the adaptor. Due to financial constraints, we are unable to evaluate larger models.

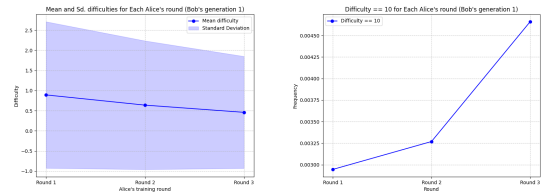


Figure 4: Instance difficulty with respect to an untrained Bob, plotted against the number of Alice’s training rounds, Qwen3-4B-Thinking-2507 models. Left: mean and standard deviation. Right: fraction of instances with maximum difficulty.

3.2 Intrinsic Evaluation

Our goal is to improve our model’s performance on code understanding tasks. In this section, we report how much better our evaluator model, Bob, performs on the semantic inequivalence game after its first and only training round. We use the final trained generator model Alice (from round 7) to generate the challenge instances. These instances are created using source programs from either the training portion of the MBPP dataset, as done during training, or from the test portion of MBPP, which neither Alice nor Bob have seen before. The results are reported in Table 1.

We observe that while the performance of the untrained Bob (baseline gpt-4o-mini-2024-07-18) is already high, this is expected because we did not train Alice to convergence. However, performing a single training round for Bob substantially improves its ability to play the game.

This demonstrates that our training protocol is effective in teaching LLMs to reason about the inputs that cause different variants of a program to

402 behave differently.

403 3.3 Qualitative analysis

404 We manually analysed the some of Alice’s outputs
405 in order to understand what kind of manipulation
406 it introduces. A typical pattern that we find is the
407 introduction of subtle bugs in conditional branches
408 that deal with edge cases. For instance, starting
409 from the following program from MBPP:

```
410 from sys import maxsize
411
412 def max_sub_array_sum(a, size):
413     max_so_far = -maxsize - 1
414     max_ending_here = 0
415     start = 0
416     end = 0
417     s = 0
418     for i in range(0, size):
419         max_ending_here += a[i]
420         if max_so_far < max_ending_here:
421             max_so_far = max_ending_here
422             start = s
423             end = i
424         if max_ending_here < 0:
425             max_ending_here = 0
426             s = i+1
427     return (end - start + 1)
```

428 Alice generates:

```
429 from sys import maxsize
430
431 def max_sub_array_sum_y(a, size):
432     max_so_far = -maxsize - 1
433     max_ending_here = 0
434     start = 0
435     end = 0
436     s = 0
437     for i in range(0, size):
438         max_ending_here += a[i]
439         if max_so_far < max_ending_here:
440             max_so_far = max_ending_here
441             start = s
442             end = i
443         if max_ending_here < 0:
444             s = i + 1
445     return end - start + 1
```

446 These programs look superficially the same,
447 they only difference is that the generated variant
448 lacks the `max_ending_here = 0` statement in the
449 `if max_ending_here < 0:` branch inside the loop,
450 which causes it to mishandle negative values in the
451 array `a`.

452 This sort of bugs often occur in programs that
453 have security vulnerabilities, therefore we believe
454 that by being trained on such examples, Bob can
455 learn to reason on the semantics of vulnerabilities.
456 This could explain the improvements we observe
457 on vulnerability datasets (Section 3.4.2).

458 3.4 Extrinsic Evaluation

459 Being proficient at playing the semantic inequiv-
460 alence game may be directly useful in certain cir-
461 cumstances, such as determining whether a code
462 refactoring has introduced subtle bugs. However,
463 ultimately, we aim for this game to teach LLMs
464 skills that generalize to a variety of tasks. There-
465 fore, we evaluate our approach across a range of
466 code-related tasks using standard benchmarks.

467 While we include code generation tasks, our
468 primary focus is on code understanding. There-
469 fore, we use the trained evaluator Bobs, denoted
470 as `sinq-gpt-4o-mini` (based on `gpt-4o-mini`),
471 and `sinq-gpt-4o-mini` (based on `gpt-4.1-nano`)
472 as our main checkpoints. Each trained model
473 is primarily compared to its own base model.
474 We do not evaluate the models based on
475 `Qwen3-4B-Thinking-2507`, as that set of experi-
476 ments yielded inconsistent training curves and was
477 aborted.

478 3.4.1 Python Builtin Identifier Swap

479 The **Python builtin identifier swap** is a very chal-
480 lenging code understanding benchmark introduced
481 by [Miceli Barone et al. \(2023\)](#). In its classification
482 version, each example consists of two variants of
483 a Python snippet, with an instruction asking the
484 model to determine which variant is more likely
485 to be correct. The challenge is that the snippets
486 are prepended with a statement that reassigns two
487 builtin Python functions used in the code, e.g.

```
488 print, len = len, print
```

489 One of the two snippets is a Python function ex-
490 tracted from a GitHub repository, while the other
491 is the same function with all instances of the two
492 builtin identifiers (e.g., `len` and `print`) swapped.
493 Because of the reassignment of the two identi-
494 fiers, the modified snippet is correct but highly
495 out-of-distribution, while the original snippet is
496 in-distribution but incorrect. [Miceli Barone et al. \(2023\)](#)
497 found that this confused all the state-of-the-
498 art LLMs available at the time, to the point that
499 they even performed worse as their size increased,
500 a case of **inverse scaling** ([McKenzie et al., 2023](#)).

501 We evaluate `gpt-4o-mini` and `gpt-4.1-nano`,
502 which had not been released at the time of the origi-
503 nal study, and our own `sinq-gpt-4o-mini` and
504 `sinq-gpt-4.1-nano` (trained Bobs) on this bench-
505 mark. We use either the original prompt template
506 or a variant that instructs the models to perform
507 chain-of-thought reasoning before answering. We

report our results in Table 2.

We observe that, despite gpt-4o-mini-2024-07-18 and gpt-4.1-nano-2025-04-14 being released years after this benchmark was published, they still performs very poorly. In fact, they perform worse than GPT-3.5 (3.35% accuracy)⁶, indicating that this benchmark remains challenging. Our approach yields a substantial improvement (+3.7%/+4.6%) over the baseline without using CoT. Surprisingly, the improvement when using CoT is smaller (+0.4%/+0.1).

This benchmark is quite different from the synthetic data used to train our model in the **semantic inequivalence game**. The main similarity is that both tasks involve reasoning about the semantics of unusual snippets of Python code. The substantial improvements we observe indicate that our approach teaches the model generalizable code reasoning skills.

We report additional results on this benchmark with state-of-the-art reasoning models in Appendix D.

3.4.2 Vulnerability Detection

Security vulnerabilities in code often arise from counterintuitive behaviours, where the intuitive understanding that programmers, whether human or LLM, have of the code’s semantics differs from its actual semantics in edge cases that evade pre-deployment testing. Our semantic inequivalence game incentivises the generator Alice to find edge cases that fool the evaluator Bob, who is then incentivised to become more robust by improving his reasoning about code semantics. Ideally, these capabilities should generalize to security vulnerability detection.

We evaluate our approach by testing it on two vulnerability detection benchmarks.

PySecDB (Sun et al., 2023) is a dataset of commits, represented as diff patches, for Python programs, classified as either containing or not containing a security fix. We present these patches to the LLMs, instructing the models to classify them. We do not provide the rest of the repository as a reference, making this a challenging task. Since some of these commits are quite long, we discard

⁶Raw results for Miceli Barone et al. (2023) are available on the GitHub repository associated with the paper: https://github.com/Avmb/inverse_scaling_prize_code_identifier_swap/blob/main/eval_chat_llms/eval_chat_llms_results.json.

those that exceed the maximum context length of 128,000 tokens.

CodeXGLUE Defect Detection (Lu et al., 2021) is a dataset of code snippets in C/C++ classified according to whether they contain known security vulnerabilities. This is a particularly challenging dataset for our approach, as we fine-tuned our model only with Python code.

We run our experiments using standard greedy classification (temperature = 0.0, no CoT), majority voting out of 9 samples (temperature = 1.0, $N = 9$, no CoT), and CoT mode (temperature = 0.6, $N = 1$). The results are reported in Table 3.

Our approach yields small but consistent improvements across two datasets, with different tasks and programming languages. These results suggest that our models have acquired additional capabilities in reasoning about security vulnerabilities, despite not having been specifically trained for this task.

3.4.3 Code Generation

We run a standard code generation experiment using the EvalPlus harness (Liu et al., 2023, 2024), which evaluates LLMs on the test portions of MBPP and HumanEval (Chen et al., 2021), as well as on augmented versions of these datasets, MBPP+ and HumanEval+, which contain additional unit tests per instance. The results are reported in Table 4.

We observe that for gpt-4o-mini our approach substantially improves code generation Pass@1 accuracy on both the original MBPP test set (+2.1%) and the more challenging MBPP+ version (+0.8%). It maintains the same level of accuracy on HumanEval and loses a slight amount of accuracy on the more difficult HumanEval+ (-0.4%). For gpt-4.1-nano we instead observe largely the same accuracy on MBPP, MBPP+ and HumanEval and a substantial improvement on HumanEval+ (+3.1%).

While our models have been trained on data from the test portion of MBPP, they have not been specifically trained to solve the MBPP task. They have never seen the natural language instructions. In fact, our models are based on the evaluators (Bobs), which have not been fine-tuned for code generation, yet they still manage to improve or maintain their generation performance. For tasks oriented towards code generation, it may be beneficial to train a separate model combining the final training datasets of both Alice and Bob.

Base LLM	Baseline	SInQ	Baseline CoT	SInQ CoT
gpt-4o-mini	1.65	5.35	1.90	2.30
gpt-4.1-nano	2.80	7.40	14.35	15.30

Table 2: Accuracy on the Python builtin identifier swap benchmark for baseline (untrained Bob) and **SInQ** (trained Bob) models, with and without chain-of-thought prompting.

Benchmark	Base LLM	Baseline	SInQ	Baseline Maj	SInQ Maj	Baseline CoT	SInQ CoT
PySecDB	gpt-4o-mini	82.43	82.51	82.48	82.81	73.55	73.00
	gpt-4.1-nano	83.20	83.33	83.03	83.35	78.74	78.20
CodeXGLUE	gpt-4o-mini	55.23	55.60	55.12	56.04	47.69	47.22
	gpt-4.1-nano	54.76	55.27	54.83	55.23	52.20	49.85

Table 3: Vulnerability detection results on the PySecDB and CodeXGLUE Defect Detection benchmarks, with or without majority voting (out of 9 samples) or chain-of-thought prompting.

Test set	Base LLM			
	gpt-4o-mini		gpt-4.1-nano	
	Baseline	SInQ	Baseline	SInQ
MBPP	82.8	84.9	87.6	86.0
MBPP+	69.6	70.4	72.5	72.2
HumanEval	87.2	87.2	89.6	90.2
HumanEval+	82.9	82.3	84.1	87.2

Table 4: Pass@1 rates (in %) on the EvalPlus suite for baseline (untrained Bob) and **SInQ** (trained Bob) models.

4 Conclusions

We presented a method to enhance the code understanding capabilities of Large Language Models by training them in a self-play setting using the **semantic inequivalence game**.

We motivated the design of this approach with theoretical arguments, demonstrating that it can cover broad domains of real-world programming by being **grounded** in a dataset of examples, while simultaneously having **no theoretical performance cap**. This allows, in principle, for unbounded performance improvements, constrained only by the available computing resources and the learning capacity of the underlying LLMs.

We evaluated our method on a variety of code reasoning tasks, including the challenging **Python builtin identifier swap** benchmark and two security vulnerability detection benchmarks. These evaluations show that our approach learns skills that generalize across tasks and programming languages.

We release all the code necessary to reproduce our experiments, along with the synthetic

training data we generated⁷. Exact replication, limited by sampling randomness, should be possible with a modest budget (approximately \$600), as long as gpt-4o-mini-2024-07-18 and gpt-4.1-nano-2025-04-14 remain available on the OpenAI platform.

We believe that our method makes a significant contribution to techniques for training LLMs on complex reasoning tasks.

Limitations

Our method has the following limitations, primarily due to our limited budget:

- We primarily fine-tuned gpt-4o-mini and gpt-4.1-nano, which, while performant, are not a state-of-the-art model. Given more resources, it would be beneficial to repeat the experiments on several more powerful models, including inference-time scaling reasoning models.
- We primarily applied supervised fine-tuning on the OpenAI platform, which likely relies on LoRA-style adaptors instead of full-parameter tuning. It would be valuable to explore reinforcement learning and full-parameter tuning as alternatives. We were able to run LoRA finetuning experiments on the small Qwen3-4B-Thinking-2507, but the results were inconsistent.
- It would be beneficial to perform multiple training rounds for Bob, with Alice being trained to convergence between each round for Bob. This could help the models learn

⁷URL redacted to preserve anonymity.

659 powerful code reasoning skills, similar to how
660 AlphaZero learns strong reasoning abilities in
661 Go or Chess through many rounds of self-play.

662 Ethical considerations

663 Our proposed method involves training LLMs on
664 synthetically generated data based on existing open-
665 source programming code datasets. We also evalu-
666 ate our method on open-source datasets. No human
667 experiments were conducted, and therefore, the risk
668 that our experiments have caused harm to individu-
669 als or infringed upon anyone’s intellectual property
670 rights is negligible.

671 Our work aims to enhance LLMs’ ability to rea-
672 son about programming code. There is a potential
673 risk that such capabilities could be used for uneth-
674 ical activities, such as hacking computer systems.
675 However, these capabilities can also be used to
676 strengthen computing systems by detecting secu-
677 rity vulnerabilities in codebases. We believe that
678 the net societal impact of our research will be posi-
679 tive.

680 References

681 Andrew W. Appel and Maia Ginsburg. 1998. *Modern*
682 *Compiler Implementation in C*. Press Syndicate of
683 the University of Cambridge.

684 Kai Arulkumaran, Antoine Cully, and Julian Togelius.
685 2019. [Alphastar: an evolutionary computation](#)
686 [perspective](#). In *Proceedings of the Genetic and*
687 *Evolutionary Computation Conference Companion*,
688 GECCO ’19, page 314–315, New York, NY, USA.
689 Association for Computing Machinery.

690 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
691 Bosma, Henryk Michalewski, David Dohan, Ellen
692 Jiang, Carrie Cai, Michael Terry, Quoc Le, and
693 Charles Sutton. 2021. [Program synthesis with large](#)
694 [language models](#). *Preprint*, arXiv:2108.07732.

695 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
696 Henrique Ponde de Oliveira Pinto, Jared Kaplan,
697 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg
698 Brockman, Alex Ray, Raul Puri, Gretchen Krueger,
699 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela
700 Mishkin, Brooke Chan, Scott Gray, and 39 others.
701 2021. [Evaluating large language models trained on](#)
702 [code](#). *Preprint*, arXiv:2107.03374.

703 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,
704 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
705 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,
706 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-
707 hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.
708 2025. [Deepseek-r1: Incentivizing reasoning capa-](#)
709 [bility in llms via reinforcement learning](#). *Preprint*,
710 arXiv:2501.12948.

Tuan Dinh, Jinman Zhao, Samson Tan, Renato Ne-
grinho, Leonard Lausen, Sheng Zha, and George
Karypis. 2023. Large language models of code fail at
completing code with potential bugs. In *Proceedings*
of the 37th International Conference on Neural In-
formation Processing Systems, NIPS ’23, Red Hook,
NY, USA. Curran Associates Inc.

Kefan Dong and Tengyu Ma. 2025. [STP: self-play](#)
[LLM theorem provers with iterative conjecturing and](#)
[proving](#). *CoRR*, abs/2502.00212.

Shihan Dou, Haoxiang Jia, Shenxi Wu, Huiyuan Zheng,
Weikang Zhou, Muling Wu, Mingxu Chai, Jessica
Fan, Caishuang Huang, Yunbo Tao, Yan Liu, Enyu
Zhou, Ming Zhang, Yuhao Zhou, Yueming Wu, Rui
Zheng, Ming Wen, Rongxiang Weng, Jingang Wang,
and 5 others. 2024. [What’s wrong with your code](#)
[generated by large language models? an extensive](#)
[study](#). *Preprint*, arXiv:2407.06153.

FAIR. 2022. [Human-level play in the game of](#)
[<i>diplomacy</i> by combining language models](#)
[with strategic reasoning](#). *Science*, 378(6624):1067–
1074.

Leonid A. Levin. 2003. [The tale of one-way functions](#).
Preprint, arXiv:cs/0012023.

Jiawei Liu, Chungju Steven Xia, Yuyao Wang, and Ling-
ming Zhang. 2023. [Is your code generated by chat-](#)
[GPT really correct? rigorous evaluation of large lan-](#)
[guage models for code generation](#). In *Thirty-seventh*
Conference on Neural Information Processing Sys-
tems.

Jiawei Liu, Songrun Xie, Junhao Wang, Yuxiang Wei,
Yifeng Ding, and Lingming Zhang. 2024. [Evaluating](#)
[language models for efficient code generation](#). In
First Conference on Language Modeling.

Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey
Svyatkovskiy, Ambrosio Blanco, Colin B. Clement,
Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong
Zhou, Linjun Shou, Long Zhou, Michele Tufano,
Ming Gong, Ming Zhou, Nan Duan, Neel Sundar-
san, and 3 others. 2021. Codexglue: A machine
learning benchmark dataset for code understanding
and generation. *CoRR*, abs/2102.04664.

Nickil Maveli, Antonio Vergari, and Shay B. Co-
hen. 2025. [What can large language models cap-](#)
[ture about code functional equivalence?](#) *Preprint*,
arXiv:2408.11081.

Ian R. McKenzie, Alexander Lyzhov, Michael Pieler,
Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan
McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, Andrew
Gritsevskiy, Daniel Wurgaft, Derik Kauffman,
Gabriel Recchia, Jiacheng Liu, Joe Cavanagh, Max
Weiss, Sicong Huang, The Floating Droid, and 8 oth-
ers. 2023. [Inverse scaling: When bigger isn’t better](#).
Preprint, arXiv:2306.09479.

Antonio Valerio Miceli Barone, Fazl Barez, Shay B.
Cohen, and Ioannis Konstas. 2023. [The larger they](#)

711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766

Theorem A.1. *There is no perfect evaluator program Bob^* such that, for any inequivalent programs P and Q it computes a diverging input for P and Q .*

Proof. If programs P and Q have a diverging input for which they both halt producing distinct output values: $P(\hat{x}) = y_p \in \mathbb{N}$, $Q(\hat{x}) = y_q \in \mathbb{N}$ and $y_p \neq y_q$, then Bob^* can compute \hat{x} by dovetailing. The interesting case is when for each diverging input only one between P and Q halts. We show that such diverging inputs cannot be computed in the general case by a reduction to the halting problem.

Given a program A and a natural number n , it is possible to algorithmically construct two programs $P_{A,n}^*$ and $Q_{A,n}^*$ defined as follows:

Listing 1: Definition of $P_{A,n}^*$ and $Q_{A,n}^*$

```

887 def P_A_n_star(x):
888     if x == 0:
889         A(n)
890         return 1
891     else:
892         while True:
893             pass
894
895 def Q_A_n_star(x):
896     if x == 0:
897         return 1
898     else:
899         P_A_n_star(x-1)

```

By construction, if A halts on input n , then $P_{A,n}^*$ and $Q_{A,n}^*$ diverge only on input 1:
 $A(n) \neq \perp \iff P_{A,n}^*(1) = \perp \wedge Q_{A,n}^*(1) = 1$,
 otherwise if A does not halt on n , then $P_{A,n}^*$ and $Q_{A,n}^*$ diverge only on input 0:
 $A(n) = \perp \iff P_{A,n}^*(0) = \perp \wedge Q_{A,n}^*(0) = 1$.
 They are always equivalent on any other input.
 Therefore:

Listing 2: Halting decider

```

908 def Halt(A, n):
909     def P_A_n_star(x): ...
910     # defined as in Listing 1
911
912     def Q_A_n_star(x): ...
913     # defined as in Listing 1
914
915     return Bob_star(P_A_n_star,
916                    Q_A_n_star) == 1

```

Since a general program that decides the halting problem cannot exist, then a perfect evaluator for the semantic inequivalence problem Bob^* cannot exist. \square

It can be noted that the proof of Theorem A.1 applies in the general case but deviates from the

constraints of the semantic inequivalence game in two important aspects:

1. The proof involves distinguishing between the halting behaviour of programs under arbitrary runtime, while in the game programs are checked against the diverging inputs produced by Alice and Bob under a time limit, after which they are assumed to return a special "TIMEOUT" value.
2. In the proof we allow both Bob's input programs P and Q to take a special form that depends on the program A whose halting behaviour is under consideration, while in the game the program P is sampled from a dataset and Alice only controls program Q .

It may be asked whether these constraints make Bob's task substantially easier, allowing for a perfect Bob^* to exist, which would imply a performance cap. We show that this is not the case.

In order to address the first point, we note that while the halting problem under a time limit is decidable if we only require the halting detector program to eventually halt, it is still undecidable if the halting detector program has to halt itself within the same time limit of the program it checks⁸. Therefore, by constraining Bob's resource usage, allowing Alice to always have more resources than Bob, and gradually increasing the time limit of the programs, it is possible for Alice to always generate harder and harder instances. Once Bob stops learning, the resource limits can be increased, enabling further learning, in principle forever. In practice, Alice and Bob are implemented as agents based on LLMs operating in chain-of-thought mode, thus resource limits can be enforced by controlling the number of reasoning tokens, or in the long term by controlling the parameter count, layer count, or expert count of the base LLMs⁹.

As for the second point, we show that for any **non-trivial program** P , Alice can generate a program $\overline{Q}_{P,A,n}$ which checks whether program A halts on input n , where by "non-trivial" we mean that there exist at least two distinct inputs x_0 and x_1 such that P halts on both, returning two distinct values, respectively y_0 and y_1 :

⁸This is provable with an argument about program length.

⁹Assuming that LLMs always become better at learning when increasing their resource limits.

Listing 3: Definition of $\overline{Q_{P,A,n}}$

```

968 def Q_P_A_n_bar(x):
969     if (x == x_0) or (x == x_1):
970         if Halt(A, n):
971             # defined as in Listing 2
972             return y_0
973         else:
974             return y_1
975     else:
976         return P(x)

```

This is a self-referential construction, where Bob is tasked to analyse a program that invokes Bob itself, thus Bob has to analyse its own behaviour.

If Bob was indeed the perfect *Bob*, then P and $\overline{Q_{P,A,n}}$ would return different values only on input x_0 if A halts on n , or only on x_1 if it does not, thus solving the halting problem. Note that this construction is still a valid output for Alice even when Bob is not perfect, since $\overline{Q_{P,A,n}}$ will still differ from P on x_0 or x_1 (possibly on both if the inner call to Bob does not halt), which means that in principle Alice can generate hard examples for Bob from arbitrary source programs, as long as they meet minimal "non-triviality" conditions. In practice, we want the generated programs to run quickly on the CPU without invoking LLMs, so this self-referential construction is unwieldy, but it serves as a proof of concept which shows that arbitrarily complex logic can be added by Alice in the programs it generates, even starting from minimally complex source programs.

B Setting a target difficulty

In the implementation of the semantic inequivalence game which we use in our experiment, we instruct the generator "Alice" to create challenge instances for the evaluator "Bob" with a specific target difficulty, defined as 10 times the probability that Bob fails to solve the instance when invoked in sampling mode. Setting the target difficulty always at the maximum value of 10 makes the game equivalent to its original formulation in section 2.1, which, if Alice never produces invalid instances, is a **zero-sum game**.

It may be asked whether this maximally adversarial setting is always ideal. Consider the following Python program that Alice may potentially generate:

Listing 4: Cryptographically hard Q generated for a given P

```

import hashlib 1014
1015
def Q(x): 1016
    try: 1017
        e = str(x).encode("utf-8") 1018
        h = hashlib.sha3_256(e).hexdigest() 1019
        t = ("af9ac3dac56b02f1ea017e7657a9bb7e" 1020
            "1778274e31509f134f023e41a5953866") 1021
        if h == t: 1022
            return "Bananas" 1023
    except: 1024
        pass 1025
    return P(x) 1026

```

For inputs x that have the specific SHA-3-256 value defined in the code, Q returns the string "Bananas", otherwise it behaves as P , therefore, as long as P does not happen to also return "Bananas" for all these specific inputs, they are diverging inputs.

Alice can easily generate this instance by first choosing a diverging input \hat{x} (in this example, the string "correct horse battery staple"), then hashing it and hardcoding its hash value into Q , but, in order to solve this instance Bob has to successfully execute a **preimage attack** on SHA-3-256, which is considered a strong cryptographic function (National Institute of Standards and Technology (NIST) and Dworkin, 2015). While this attack is theoretically possible by brute-force search, in practice it would require a runtime longer than the age of the universe, unless perhaps Bob is a cryptanalysis genius and manages to find a serious flaw in SHA-3-256, and even in this case, if the **one-way function conjecture** happens to be true then it is possible to construct asymptotically strong cryptographic hash functions (Levin, 2003), making Bob's task effectively hopeless.

The construction used is our specific example would require Alice to run code in order to compute the hash of its chosen diverging input, which current LLMs are typically not allowed to do in their default configuration and not in our experiments (although some common "LLM agent" setups do allow it), but Alice could still manage to create cryptographic puzzles which are too hard for any practical Bob to solve.

If Alice is instructed to always generate maximally difficult instances, it has an incentive to generate cryptographic puzzles, but since Bob only learns from the instances it can actually solve, this would effectively cause the learning process to stall. In Appendix A we have proven that learning can continue forever in the limit of infinite computing

resources, but in reality computing resources are finite, and cryptographic puzzles could stop the learning process as soon as Alice discovers the trick. Even if it never resorts to cryptographic puzzles, Alice could just learn faster than Bob, eventually overwhelming Bob with instances that it cannot solve and thus stopping the learning process.

Fortunately, we can avoid this problem completely by setting the target difficulty to a value lower than the maximum, e.g. 7, corresponding to the current Bob solving the instances with 30% probability. This changes the nature of the game from **zero-sum** to **positive-sum**, where Alice acts as a teacher that challenges Bob with instances which are hard, but not too hard for its current level. As Bob improves, the difficulty of a given distribution of instances decreases, which in turns causes Alice to learn to recalibrate its difficulty estimation and gradually generate more challenging instances, enabling the training process to continue learning interesting coding logic for as long as the capacity of the underlying LLMs is not exceeded.

In our experiments, due to our limited resources, we could not train Alice to the point that it could seriously challenge Bob, thus we always set the target difficulty to 10, but as a training recipe, we do recommend reducing the target difficulty if at some point Bob starts to fall behind.

C Prompt templates

System prompt for Alice

You are an expert computer scientist. Your task is to take a Python 3.10 program and write a similar program which is not semantically equivalent, which means that there must exist at least a diverging input example such that the original program and your program either produce different outputs or exceptions, or one halts and the other one does not halt. In addition to a program, you need to produce a diverging input example. Start by carefully analyzing the original program and think of how an example would propagate through it from the input to the return value, considering how to modify the program in order to elicit a different behavior. Make sure that the return values or exceptions raised by your program are picklable. The original program and your program will be used in a test to evaluate the skill of an expert computer scientist who will have to produce a diverging example (not necessarily the same as yours), so make sure that the difference you introduce are not very easy to understand. You will be given a difficulty level from 0 (easiest) to 10 (hardest) to target. E.g. difficulty level 0 means that an expert computer scientist in the bottom decile or above should be able to find a diverging example, difficulty level 9 means that only an expert computer scientist in the top decile should be able to find a diverging example, and

difficulty level 10 means that only the top 1% or less of expert computer scientists should be able to find a diverging example. Think step by step before writing your program. Use the following Markdown format, making sure that the following sections are delimited by level 1 headings, since they will have to be automatically parsed:

```
# Analysis
step by step analysis. This section can include
sub-headings and code blocks
# Generated program
your program inside a Python code block. Do not
change the name or signature of the entry point
function
# Diverging input example
your diverging input example as a Python dictionary
inside a Python code block
For instance, if the entry point function takes
two parameters a and b and your diverging example
is a="foo" and b=42, write:
```python
{
 "a": "foo",
 "b": 42
}
```
do not write the expected outputs
```

| | | |
|------|---|------|
| 1153 | User message for Alice As a Python f-string: | 1214 |
| 1154 | <code>f"""Difficulty level: {difficulty_level}</code> | 1215 |
| 1155 | <code>Entry point function: {function_name}</code> | 1216 |
| 1156 | | 1217 |
| 1157 | <code>```python</code> | 1218 |
| 1158 | <code>{code}</code> | 1219 |
| 1159 | <code>```"""</code> | 1220 |
| 1160 | During inference <code>difficulty_level</code> is the tar- | 1221 |
| 1161 | get difficulty (always 10), during SFT training, for | |
| 1162 | Alice’s main examples it is the measured difficulty | |
| 1163 | approximated to the nearest integer, for Alice’s dif- | |
| 1164 | ficulty prediction examples it is the string "Any". | |
| 1165 | Second user message for Alice Used only for | 1222 |
| 1166 | the difficulty prediction training examples. | 1223 |
| 1167 | Predict the difficulty level of the instance. | 1224 |
| 1168 | Just write "Difficulty level: D" where D is your | 1225 |
| 1169 | prediction, do not write anything else. | 1226 |
| 1170 | Second assistant message for Alice Used only | 1227 |
| 1171 | for the difficulty prediction training examples. As | 1228 |
| 1172 | a Python f-string: | 1229 |
| 1173 | <code>f"""Difficulty level: {difficulty_level}"""</code> | 1230 |
| 1174 | where <code>difficulty_level</code> is the measured diffi- | 1231 |
| 1175 | culty. | 1232 |
| 1176 | System prompt for Bob | 1233 |
| 1177 | You are an expert computer scientist. Your task | |
| 1178 | is to take two Python 3.10 programs and determine | |
| 1179 | whether or not they are semantically equivalent. | |
| 1180 | Two programs are semantically equivalent if there | |
| 1181 | exists no diverging input examples such that the | |
| 1182 | original program and your program either produce | |
| 1183 | different outputs or exceptions, or one halts and | |
| 1184 | the other one does not halt. If you determine that | |
| 1185 | the two programs are not semantically equivalent, | |
| 1186 | you also need to produce a diverging input example. | |
| 1187 | Start by carefully analyzing the two programs | |
| 1188 | and think of how an example would propagate | |
| 1189 | through them from the input to the return value, | |
| 1190 | considering whether it could elicit a different | |
| 1191 | behaviors. | |
| 1192 | Think step by step before writing your program. | |
| 1193 | Use the following Markdown format, making sure that | |
| 1194 | the following sections are delimited by level 1 | |
| 1195 | headings, since they will have to be automatically | |
| 1196 | parsed: | |
| 1197 | <code># Analysis</code> | |
| 1198 | step by step analysis. This section can include | |
| 1199 | sub-headings and code blocks | |
| 1200 | <code># Equivalent?</code> | |
| 1201 | Yes or No | |
| 1202 | <code># Diverging input example</code> | |
| 1203 | your diverging input example as a Python dictionary | |
| 1204 | inside a Python code block, or nothing if the two | |
| 1205 | programs are equivalent. | |
| 1206 | For instance, if the entry point function takes | |
| 1207 | two parameters <code>a</code> and <code>b</code> and your diverging example | |
| 1208 | is <code>a="foo"</code> and <code>b=42</code> , write: | |
| 1209 | <code>```python</code> | |
| 1210 | <code>{</code> | |
| 1211 | <code> "a": "foo",</code> | |
| 1212 | <code> "b": 42</code> | |
| 1213 | <code>}</code> | |
| | <code>```</code> | |
| | do not write the expected outputs | |
| | Note that we ask Bob to determine whether the | |
| | two programs are equivalent, even though they | |
| | never are. This is not strictly necessary, but it poten- | |
| | tially makes the task slightly more difficult for Bob, | |
| | which is beneficial since Bob tends to be much | |
| | stronger than Alice. | |
| | User message for Bob As a Python f-string: | 1222 |
| | <code>f"""Entry point function: {function_name}</code> | 1223 |
| | | 1224 |
| | Program 1: | 1225 |
| | <code>```python</code> | 1226 |
| | <code>{code_P}</code> | 1227 |
| | <code>```</code> | 1228 |
| | Program 2: | 1229 |
| | <code>```python</code> | 1230 |
| | <code>{code_Q}</code> | 1231 |
| | <code>```"""</code> | 1232 |
| | | 1233 |
| | The evaluation prompts will be included in the | 1234 |
| | code released upon publication. | 1235 |
| | D Additional Python builtin identifier | 1236 |
| | swap results | 1237 |
| | Main results Same as Table 2, presented as a bar | 1238 |
| | chart in Figure 5. | 1239 |
| | Results on Reasoning Models Large Reasoning | 1240 |
| | Models (LRMs) are LLMs which have been specif- | 1241 |
| | ically trained to solve reasoning tasks, primarily in | 1242 |
| | the domains of math and coding, using Chain-of- | 1243 |
| | Thought reasoning. These models, such as OpenAI | 1244 |
| | <code>o1</code> and <code>o3</code> and DeepSeek-r1 (DeepSeek-AI et al., | 1245 |
| | 2025) typically generate a large amount of reason- | 1246 |
| | ing tokens during inference, hence they are said | 1247 |
| | to perform inference-time scaling by trading off | 1248 |
| | speed and cost for quality. In practice, they are very | 1249 |
| | strong but also very expensive. | 1250 |
| | Our approach could be broadly considered a type | 1251 |
| | of LRM, since it is trained to solve reasoning prob- | 1252 |
| | lems using CoT, though in practice we use a much | 1253 |
| | smaller base LLM and we do not invest nearly as | 1254 |
| | many resources neither during training time nor | 1255 |
| | during inference time. | 1256 |
| | We evaluate the OpenAI LRMs <code>o1-2024-12-17</code> | 1257 |
| | and <code>o3-mini-2025-01-31</code> and DeepSeek LRMs | 1258 |
| | <code>r1</code> and its distilled version based on Meta | 1259 |
| | <code>Llama-3.3-70B-Instruct</code> on the Python builtin | 1260 |
| | identifier swap benchmark. Due to the high cost | 1261 |
| | and low speed of inference, for all these models ex- | 1262 |
| | cept <code>o3-mini-2025-01-31</code> we only evaluate 10% | 1263 |
| | of the test set. For the OpenAI models we evalu- | 1264 |
| | ate both on the default prompt and the CoT-style | 1265 |

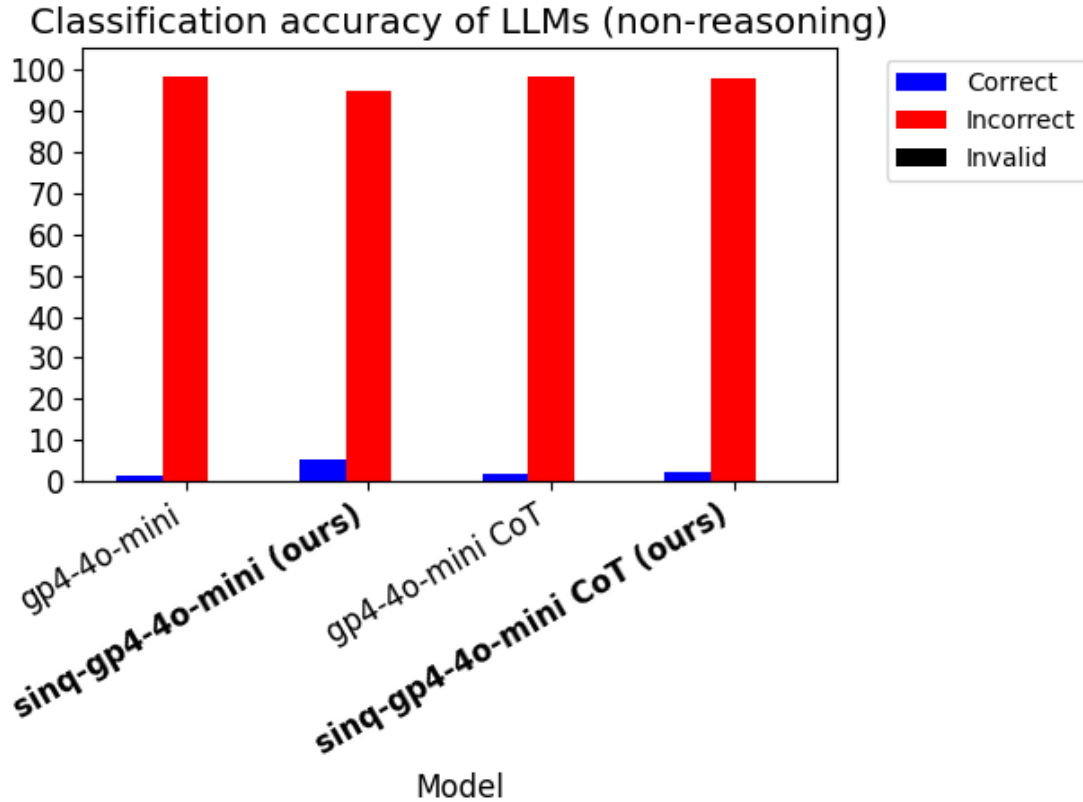


Figure 5: Python builtin identifier swap results for the baseline gpt-4o-mini (untrained Bob) and our model sinq-gpt-4o-mini (trained Bob), with or without chain-of-thought.

1266 prompt suggested for DeepSeek-r1. We report the
 1267 results in Figure 6.

1268 The LRMs are much stronger than
 1269 gpt-4o-mini, gpt-4.1-nano and our ap-
 1270 proach, with the full DeepSeek-r1 reaching 94.0%
 1271 accuracy, which is expected given their training
 1272 and inference costs.

1273 Given sufficient resources, it would be beneficial
 1274 as a future experiment to use one of these models
 1275 as the base model for our approach. We expect
 1276 that our approach would be complementary to the
 1277 synthetic data generation techniques used to train
 1278 these models, resulting in further improvements.

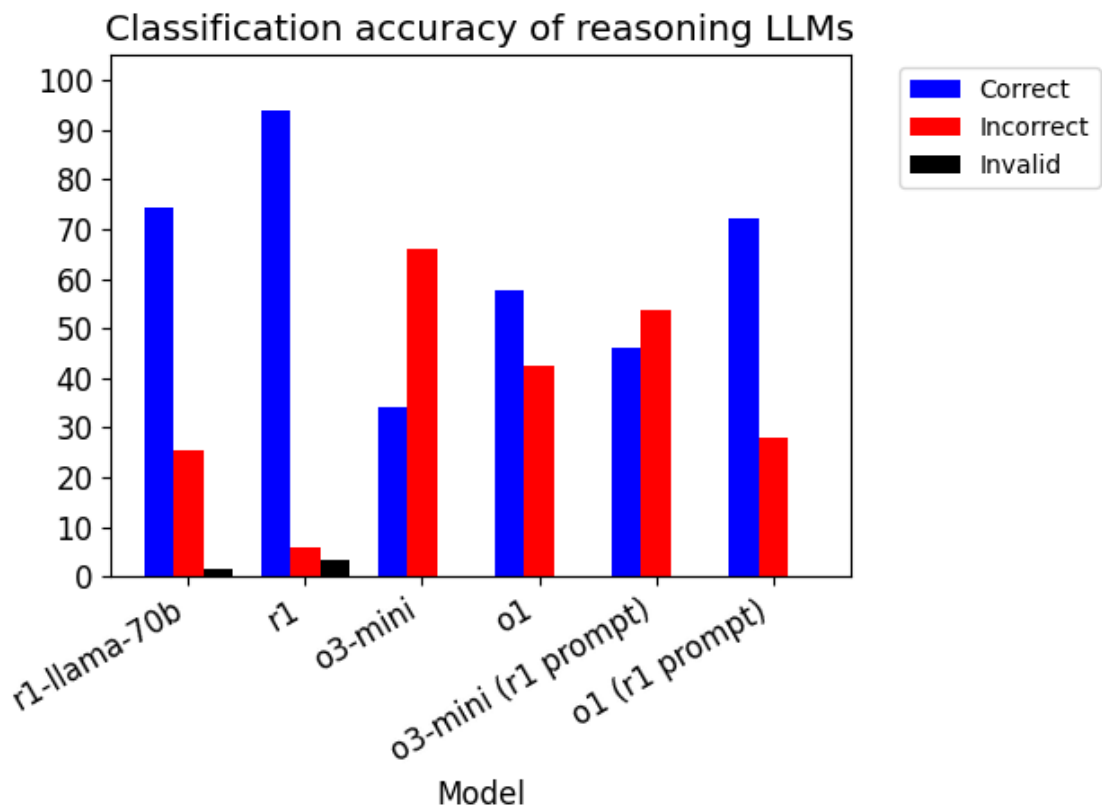


Figure 6: Python builtin identifier swap results for LRMs.