

# Inferring Code Correctness from Specification

Florian Tambon  
florian.tambon@uni.lu  
University of Luxembourg  
Luxembourg, Luxembourg

Mike Papadakis  
michail.papadakis@uni.lu  
University of Luxembourg  
Luxembourg, Luxembourg

## Abstract

Large language models (LLMs) have become integral to modern software development, enabling automated code generation at scale. However, validating the correctness of LLM-generated code remains a critical and largely unsolved challenge. Existing approaches either rely on dynamic consensus across multiple code candidates - making them costly and difficult to scale - or on static reasoning that is susceptible to dynamic bugs and order bias. In this paper, we propose TRAILS (Targeted Reasoning Agreement via Inputs and Specifications), an approach that grounds LLM reasoning with concrete (input, output) pairs. TRAILS first generates diverse test inputs via category partitioning based on the specification, then executes them against the candidate code and prompts LLMs to assess whether the resulting input-output pairs conform to the specification - without ever reasoning over the code itself. Scores are aggregated across inputs, to determine whether the program is likely correct. We evaluate TRAILS on two datasets, LiveCodeBench and CoCoClaNeL, across three LLMs (Qwen3Coder-30B, Devstral-Small-24B, and Olmo3.1-Instruct), comparing against HoarePrompt and a Zero-Shot Chain-of-Thought baseline. TRAILS improves Matthew Correlation Coefficient by up to 39% relative to Zero-Shot COT and consistently outperforms HoarePrompt. Beyond accuracy, TRAILS demonstrates greater stability across seeded runs, reducing sensitivity to LLM non-determinism, and assigns correct labels to a larger set of unique code samples than competing approaches.

## Keywords

Automated code generation, LLMs, Prompt engineering, Code correctness

### ACM Reference Format:

Florian Tambon and Mike Papadakis. 2018. Inferring Code Correctness from Specification. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

## 1 Introduction

Large language models (LLMs) have revolutionized software development by automating code generation at an unprecedented scale. Tools like GitHub Copilot, ChatGPT, and Claude now assist developers in generating code snippets, functions, and even entire

classes. Recent surveys indicate that LLM adoption in code generation workflows has grown exponentially, with industry-wide deployment across startups and enterprises alike [7]. This capability has transformed the landscape of software engineering, shifting the focus from code writing to code review and validation.

However, the widespread adoption of LLM-based code generation has revealed a critical gap: while these models excel at generating plausible-looking code, validating its correctness with regard to the original specification remains error-prone, complex, expensive and time-consuming. For experienced developers traditionally debugging and testing is already a complex and time consuming task [33], particularly when working with unfamiliar problem domains or complex specifications. With emergence of LLMs, developers starting using LLMs to offset this workload, resulting in critical step being skipped such as testing and verifying generated code [11]. In parallel, a new category of (non-expert) developers with little to no experience emerged, empowered by the capability of LLMs. However, these users face added challenges compared to developers with actual experience and knowledge, and tend to overly on LLMs given their lack of expertise as well as overtrust LLMs' answers [11–13]. The fundamental challenge underlying this gap is linked to the test oracle problem, with the absence of reliable oracle for LLMs generated code. In traditional software development, engineers invest significant effort in creating test suites before or alongside code development. While LLMs can generate test suite [8, 15, 30, 38], they are prone to mistake, even more so when the underlying code is incorrect [17, 19]. This creates a compound effect where incorrect code can be assessed as correct by incorrect test.

Given this issue, recent researches have studied the possibility to verify generated code correctness. All these studies have the same key element, which is to establish whether an inconsistency exists between the code and specification, either dynamically or statically. One of the first approach, CodeT [6], calculate a consensus set between multiple generated test cases and candidates test programs, selecting a likely correct code based on the highest consensus. In a similar fashion, Valentin et al. [35] used multiple candidates programs over a large set of generated inputs to calculate incoherence between candidates solution in order to select a likely correct candidate or abstain. Both approach relies on concrete executions of the code, and so are reliant on the generated inputs. Moreover, the approach can be hard to scale, as it relies on multiple test candidates for the same specification. Finally, their underlying hypothesis is that LLMs' incorrect codes will tend to have a wide enough array of incorrect behavior, so that incoherence with the set of correct codes will be noticeable. On the other hand to this approach, HoarePrompt [5] aims to leverage LLMs reasoning guided via Hoare's logic in order to statically verify the correctness of a code, by checking at each step the state of the program for inconsistency. While more flexible than previous approach, because

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*Conference acronym 'XX, Woodstock, NY*

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/2018/06  
<https://doi.org/XXXXXXXX.XXXXXXX>

of the absence of concrete execution, it can fail to identify dynamic bugs. Moreover, because of the inherent single binary final decision, models can be prone to order bias [31] which can hamper stability of the approach.

To tackle this issue, we propose TRAILS (Targeted Reasoning Agreement via Inputs and Specifications), which uses concrete input execution results to ground the LLMs reasoning. TRAILS follows a dual scheme in which test inputs are first generated using category partitioning based on the specification to force diversity and trigger incorrect program behavior. Inputs are then validated by running them against the candidate code, to ensure that they conform to the program input format. Then, in the second step, test inputs are executed on the program to provide candidate output (program behavior). LLMs are then prompted to reason over the specification against the input-output pairs and characterize them as conforming or not. This step pushes the LLM to reason over the program behavior with respect to the specification and identify incorrect outputs.

The advantage of our approach is that it forces the reasoning to occur exclusively on the specifications and not on the generated code, which may be incorrect and bias the LLMs responses [17, 19]. The results of multiple inputs are then aggregated, yielding the expectation that the program is correct if LLMs consider its outputs to be conforming to the specifications more often than not. As such, we leverage a threshold to determine our confidence on correctness, i.e., above a specific number of non-conforming input-output pairs the code is considered as likely incorrect.

In a way, our approach tackles the problem as a "low trust" system, where the only source of truth is emphasized to be the specification, and each subsequent step is treated with some guarding measures. To this end, we study the following Research Questions (RQs):

- RQ1** How does TRAILS 's performance compare to current approaches?
- RQ2** What is the cost of each step of TRAILS?
- RQ3** How stable is the decision of TRAILS across reruns?
- RQ4** How many consistently and correctly labeled code do each approach label?

We evaluate our approach on two different code datasets, LiveCodeBench [26] and CoCoClaNeL [5] on three different LLMs: Qwen3Coder-30B [34], Devstral-Small-24B [1] and Olmo3.1-Instruct [27], all relatively recent and capable open-source models that can git on a consumer grade single GPU. We compare our approach against current state-of-the-art approaches on correctness inference of a single sample, HoarePrompt, as well as a baseline using Zero-shot Chain-of-Thoughts (COT) Reasoning [37]. Our results show that TRAILS improves performance in terms of Matthew Correlation Coefficient across the board and up to 39% relatively to Zero-Shot COT and improves over HoarePrompt. Token consumptions mainly stems from additional repair on either intractable incorrect code (by design to ensure validity of input) or invalid inputs on CoCoClaNeL which expects standard input (stdin) for the code which models struggle with. TRAILS not only improves performance, but also is more correctly stable: over seeded runs, it preserves more the label assigned to a code than competing approaches. This makes our approach more robust to non-determinism of LLMs and limit the

effect of the seed. Finally, we show that TRAILS consistently assign the correct labels to more unique code across datasets and models compared to baselines. There is nevertheless a sizeable overlap which can motivate further research on how to properly combined different approaches.

## 2 Approach

### 2.1 Motivating example

Figure 1 illustrates how TRAILS work on an example from the CoCoClaNeL dataset. The task is to code a program that determines who between Bob and Alice will win a coin game. We compare our approach with traditional Zero-Shot COT. For traditional Zero-Shot COT, we provide the LLM with specification + the solution and ask the LLM to say whether the solution is correct or not with regard to the specification. The model incorrectly judge the code as correct. Even if it resorts to examples, it has no grounding output and so the LLM tends to follow on what the code does rather than reasoning over the specification as instructed to. In that case, it declares "Bob" to be the output for the input "a = 2, b = 1", even though the specification highlights that Alice goes first and so she should win. This behavior, LLM focusing on actual behavior of code rather than expected by specification, was already noted by previous research on test oracles [19]. On the contrary, using our proposed TRAILS, the LLM is provided with the specification, an input generated at an earlier stage and the executed output over the code solution. In that case, the LLM is asked to verify that the input-output pair is correct in relation to the specification. Since there is no code, it has to reason over the specification with the given input and, given the grounding output, highlight the contradiction between the expected behavior and actual one. In that case, it correctly exposes the bug.

### 2.2 Problem formulation

We consider the following problem. Given a specification in natural language (docstring, plain text etc.) and a code generated by LLMs which correctness is unknown, the goal is to infer the correctness of the code. Importantly, we consider that we have access to *no* oracle, that is no available test case, no example nor user feedback to draw information from. This represent a general case of a user obtaining a code with little way of ensuring validity. As described in introduction, contrary to general benchmark such as HumanEval [23] which are composed of small code snippet where some test cases could easily be derived, practical usage of code LLMs often lead to code snippet where the user can not easily infer a proper oracle (especially, non-expert). While one could use possible test cases provided by the LLMs, studies [17, 30] showed that this prone to error, especially when incorrect code is involved. Moreover, even with correct code, LLMs tend to generate incorrect or incomplete assertion, decreasing their accuracy, as noted in HoarePrompt study [5]. Essentially, we set ourselves in a worst-case/low trust settings where the only source of trust is the specification. Moreover, contrary to current dynamic approaches which relies on the availability of  $N$  code snippets, we only consider a *single* code snippet, as asking for several code snippet is not easily scalable to larger code snippet.

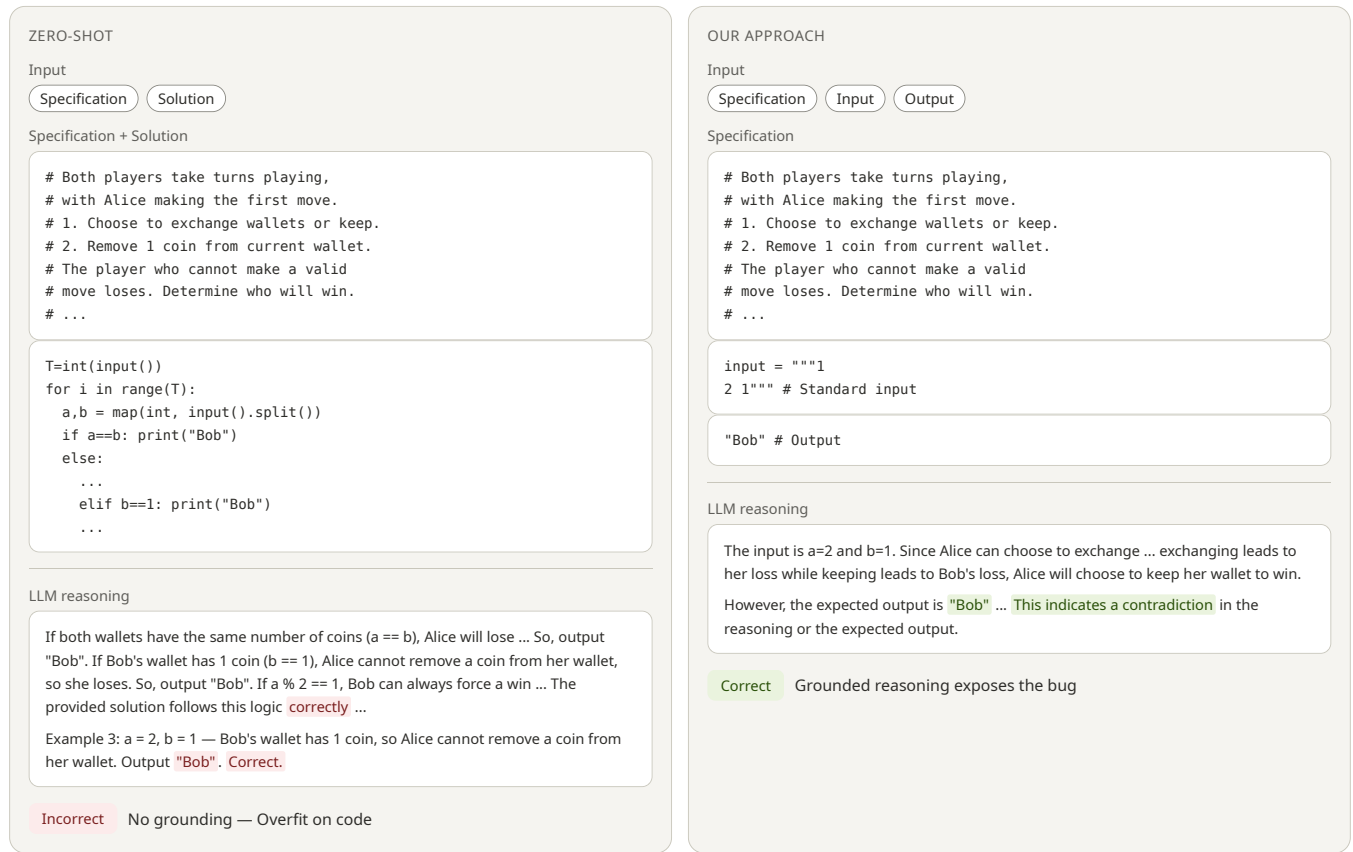


Figure 1: Motivation example. (Left) Zero-Shot reasoning: the model incorrectly validate the code logic. (Right) TRAILS correctly describe that the (input, output) pair can not be obtained with the given specification.

### 2.3 Overview of the approach

An overview of the approach is represented on Figure 2. The approach is divided into two phases: 1) input generations and 2) output verification. First, TRAILS aims to generate inputs that will be used to ground the reasoning of LLMs. To do so, it relies on categorical partitioning to extract relevant behaviors and pre-conditions from the specification. The partitions are then used to obtain inputs encompassing different behavior paths. Given inputs generated might not be valid, we adopt a simple repair/validation mechanism following existing studies on test generation [20]. Essentially, the inputs are repaired under a given budget and deemed valid if they pass on the code under test, otherwise they are discarded. In the second stage, collected inputs are run on the concrete code and outputs are collected. Then, an LLM is prompted to verify the triplet (input, output, specification) via a binary decision. This verification is done for every inputs generated, resulting in an expected probability of the code being correct. Intuitively, we assume correct codes (i.e. code matching the specification) to lead to more verified outputs than incorrect codes.

### 2.4 Framing the task inputs

The aim of this step is to obtain an array of inputs that encapsulate the distribution of possible inputs the task is requiring. Asking

directly an LLM to generate inputs might lead to edge cases being missed as the LLM will focus on common cases [38]. These edge cases might exactly be why the code of unknown correctness is not matching the specification. Applying mutations on top of generated inputs similarly to HumanEval+ [23] to increase diversity is not desirable in our case as it implies writing mutators for a large number of object types and use cases as well as verifying a high number of generated inputs with potentially broken conditions due to the mutators not being able to handle all possibilities. Given our aim is not develop a comprehensive set of inputs, but only to find possible discrepancy between code and specification, we rely on this much cheaper alternative.

As such, we first prompt an LLM to map the possible behavior scenario based on the provided specification (1). The LLM is tasked with extracting a behavior, encompassing possible pre-conditions (e.g. "n should be strictly positive") as a free form text. Contrary to using logic formula, the free form text allows the LLM to represent more complex cases which might not be easily represented using such abstraction (e.g. handled exception in the code which should raise a dedicated error). In parallel, we also query an LLM to extract relevant properties of the code under test (1). This step is important to be able to properly execute the inputs down the line, providing the model with all relevant information. In particular,

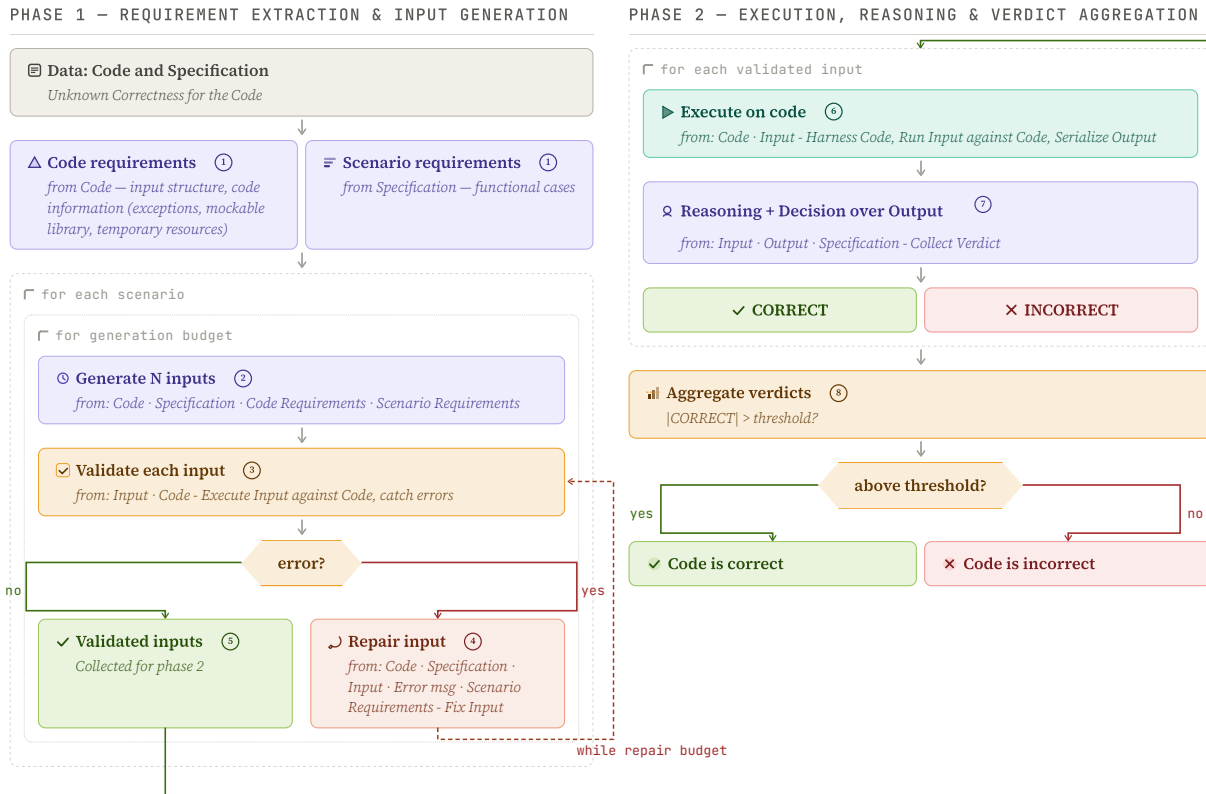


Figure 2: TRAILS overview

besides simple use case where inputs is either a list of arguments or a structure stdin, inputs could also require mockable dependency, temporary files, exception handling etc. For instance, when tackling a function sending GET request, it is imperative to be able to control the effect of the request to be able to cover all possibilities of the task, which is not doable via simple input injection. As such, the LLM is tasked with extracting properties such as the inputs structure, possible exceptions handling, mockable object or temporary resources needed.

## 2.5 Generation of valid input inputs

Then, given an extraction partition behavior, the LLM is prompted to generate proper inputs for the code under test while accounting for the particular properties extracted (2). Then, each generated input is then validate by running it against the code under test (3). Any input leading to an unhandled crash of the code is deemed invalid. While this might exclude potentially valid and relevant inputs, this is nonetheless necessary following our low-trust philosophy: if a crash occurs, it might be because of a flawed code or an issue with the input, which is not trivial to decide. Nonetheless, given LLM might struggle on generating a valid input on some tasks, we add a simple repair mechanism with a given budget following current practice on test repair [20]. Essentially, the LLM is prompted with the current input and the error raised and asked to correct the input, assuming a correct program (4). If no valid input can be generated in a given budget, even with repair, the partition is skipped and

a new one is processed. If consecutive partitions lead to no valid input, an early stopping mechanism declare the code under test as invalid, preventing additional computation on a likely incorrect code. This steps leads, for each partition  $p$  obtained, to a number  $I_{tot,p}$  of inputs. To further optimize the process, a deduplication step (not shown in the graph) is ran by computing, within each partition, code coverage to remove inputs triggering the same code lines. As such, one obtain a number of inputs per partition  $I_{reduced,p} < I_{tot,p}$  (5).

## 2.6 Output verifications

The code is harnessed and the obtained inputs  $I_{reduced,p}$  are ran against it to obtain the outputs  $O_{reduced,p}$  (6). Each output is mapped onto an object that can be easily serializable for the verification step.

Given the triplets  $(I_{reduced,p}, O_{reduced,p}, \text{specification})$ , an LLM is prompted using COT to verify the output using binary labels outcome ("CORRECT"/"INCORRECT") and to explain its reasoning (7). While generating directly the oracle  $O$  via reasoning is doable and could be potentially more informative, multiple studies highlight the difficulty for LLMs to generate proper oracle [5, 17, 38], even more so when complex types are expected. Comparatively, using verification grounded on a concrete (input, output) pairs simplifies the process for the LLM. Intuitively, this also provide a grounding reference compared to plain Zero-Shot COT. Moreover, given we

have multiple inputs per partition, we offset possible stochasticity issue due to binary decision.

## 2.7 Scoring over inputs

In the end, the verdicts are aggregated we obtain a score per code under test (8). The obtained score represents the percentage of agreements across the inputs and is an approximation of the expected probability that the code under test is correct. Intuitively, the more agreement over inputs, the more likely the code is indeed correct. This aggregated score is compared to a threshold to decide whether the code is correct or not. This threshold might be obtained in two ways. The most straight forward is to use establish threshold over benchmarks. Since the score is in  $[0, 1]$ , a higher threshold will lead to a more conservative approach (harder to accept a code as correct) at the cost of potentially discarding valid code. This is a trade-off that can be tuned by the user. An alternative is to tune the threshold on a calibration set that is provided based on the task to assess. While this might be harder to implement (it requires a dedicated calibration set), it gives more accuracy in the choice of threshold. In our experiment, we picked the threshold empirically, as experiments showed that one threshold gave best performance across models and datasets.

## 3 Experimental setup

### 3.1 Datasets

In our study, we use two datasets. LiveCodeBench [26] is a dataset of competition programming extracted from different website, containing problem description with hidden test cases but no ground truth implementation. In this study, we use the v4 Lite versions, resulting in 119 tasks from half of 2024. To generate code implementations, we used a QwenCoder2.5-7B [18] which was prompted to solve each task following the prompts given in the respective benchmarks. The model was prompted 10 times and obtained code where executed against the tests available in the benchmarks. For each task, we recovered a correct sample or an incorrect sample (or both) depending on the model performance. This model was not used in our study and is older than the LLMs used for evaluation to prevent bias. We also included CoCoClaNeL [5] for comparison as it was introduced by HoarePrompt. It includes 163 tasks with 654 problem-solution pairs from Codeforces programming contests that took place in the first half of 2024, with annotated ground truth. Contrary to LiveCodeBench, this benchmark include human code of varying quality. Manually analyzing the obtained set, we noted 2 tasks with incomplete problem description which were removed from the set resulting in 161 tasks. For this dataset, we sampled two code (one correct/one incorrect) for each task, when possible (some tasks have no correct samples).

### 3.2 Models

To evaluate our approach, we chose four models encompassing different architecture, training process and providers. As some baseline use their own reasoning as part of the process, to draw a fair comparison, we chose only models with no reasoning capability (or with capability turned off) to avoid influence of the reasoning over the results. We chose Qwen3Coder-30B [34], Devstral-Small-24B [1] and Olmo3.1-32B-Instruct [27]. The motivation to chose these particular

models are multiple: they are all are mid-size open-source LLMs capable on fitting on a single consumer-grade GPU, they are from different providers (hence different architecture, training pipeline etc.) and do not have high performance on LiveCodeBench (and so CoCoClaNeL) contrary to frontier closed-source models which allow to evaluate the usefulness of the approach rather than absolute models strength. Temperature was fixed to 0.5 for all models following HoarePrompt default settings.

### 3.3 Approaches

As a baseline, we used both Zero-Shot Chain-of-Thoughts (COT) [37] and state-of-the-art HoarePrompt [5]. For Zero-Shot COT, we use traditional prompt with the model needing to first reason (e.g. "Think step-by-step") before giving a verdict between "CORRECT" and "INCORRECT". The model is only given the code under test (of unknown correctness) and the specification, and must decided whether the model implement the specification. For HoarePrompt, we follow default parameters which use the whole pipeline (code + spec, no unrolling, precondition extraction, natural strongest postcondition and correctness classification) and also return "CORRECT" or "INCORRECT". We do not make any change to the prompts used. For TRAILS, we set the number of scenario to be generated to  $s = 3$ , the generation budget to  $g = 3$  (i.e. at most 3 input/repair steps) and maximum deduplicated inputs per scenario to  $|I_{reduced,p}| = 3$ . In all cases, the temperature is set to 0.5 following HoarePrompt default configuration.

## 4 Results

### 4.1 RQ1: How does TRAILS 's performance compare to current approaches?

In our first research question, we compare the performance of TRAILS against other baselines. To track the performance, we use both Matthew Correlation Coefficients (MCC) [39] similarly to HoarePrompt's study, as well as P4 metric [32]. P4 is a symmetric version of the F1 metric and considers also the whole confusion matrix as MCC does. Contrary to MCC, P4 put more emphasis on extreme cases (e.g. few false positive etc.) which is useful to contrast MCC value. Results are averaged over three reruns across datasets for each method/model. Results are given in Table 1.

Overall, TRAILS improves performance over all models and datasets, yielding more correct classifications. We compared thresholds and display three best in terms of metrics.  $\tau = 0.8$  seems to be the best across the board with  $\tau = 0.7$  being a close second.  $\tau = 0.6$  have degraded performance compared to Zero-shot especially in terms of P4. At similar model/dataset, relative performance gain is more emphasize on CoCoClaNeL compared to the Zero-shot COT. This might simply be explained by LiveCodeBench's task being more simple, and so there is less of an improvement to be made, the ceiling of performance being the performance of the model itself. We note the degraded performance of HoarePrompt on LiveCodeBench. We explain this decrease again by the "simpler" tasks of LiveCodeBench: in that case, the added logical reasoning might artificially complexify the reasoning trace compared to simple CoT, degrading performance. On the contrary, CoCoClaNeL involves coding competition code, written by human and so more complex specification as well as harder to process

**Table 1: Comparison of LLMs across two datasets using our approach and baselines. Best result per dataset–metric block is in bold; second-best is underlined. Results are averaged over three reruns. Tokens are given in thousands per task. Increase in performance is given relative to Zero-shot CoT.**

Dataset	Model	TRAILS									HoarePrompt			Zero-shot CoT		
		$\tau=0.6$			$\tau=0.7$			$\tau=0.8$			MCC (†)	P4 (†)	Tokens (‡)	MCC (†)	P4 (†)	Tokens (‡)
		MCC (†)	P4 (†)	Tokens (‡)	MCC (†)	P4 (†)	Tokens (‡)	MCC (†)	P4 (†)	Tokens (‡)						
LCB	Qwen3-Coder	.633 + 3.43%	.776 - 3.48%	22,1k	<u>.655</u> + 7.03%	<u>.806</u> + 0.24%	22,1k	<b>.661</b> + 8.01%	<b>.825</b> + 2.61%	22,1k	.605 - 1.14%	.798 - 0.06%	11,8k	.612	.804	1,1k
	Devstral-Small2	<b>.550</b> + 18.79%	.737 + 3.37%	18,6k	<u>.536</u> + 15.76%	<u>.753</u> + 5.61%	18,6k	.508 + 9.71%	<b>.754</b> + 5.75%	18,6k	.357 - 22.89%	.660 - 7.43%	16,7k	.463	.713	1,2k
	Olmo3.1-Instr	.579 + 24.78%	.773 + 6.92%	21,9k	<u>.601</u> + 29.52%	<u>.793</u> + 9.68%	21,9k	<b>.606</b> + 30.60%	<b>.803</b> + 11.07%	21,9k	.355 - 23.49%	.678 - 6.22%	11,3k	.464	.723	1,3k
CoCo	Qwen3-Coder	.211 + 13.44%	.500 - 2.34%	37,5k	<u>.259</u> + 39.24%	<u>.591</u> + 15.43%	37,5k	<u>.246</u> + 32.26%	<b>.617</b> + 20.51%	37,5k	.223 + 19.89%	<u>.608</u> + 18.75%	15,0k	.186	.512	1,7k
	Devstral-Small2	.214 + 0.94%	.521 - 11.99%	27,0k	<u>.247</u> + 16.51%	<u>.595</u> + 0.51%	27,0k	<b>.261</b> + 23.11%	<b>.630</b> + 6.42%	27,0k	.223 + 5.19%	<u>.609</u> + 2.87%	20,6k	.212	.592	1,5k
	Olmo3.1-Instr	.317 - 4.51%	.598 + 1.18%	30,3k	<u>.411</u> + 23.80%	<u>.684</u> + 15.74%	30,3k	<b>.431</b> + 29.82%	<b>.711</b> + 20.30%	30,3k	.269 - 18.98%	.627 + 6.09%	14,1k	.332	.591	2,5k

code. For instance, such code snippets often contains meaningless variable identifier ("a = ..., ii = ...") which were shown to degrade models' performance in reasoning [36]. In that instance, TRAILS's triplet (input, output, specification) can help mitigate the reasoning degradation by grounding it with concrete inputs. Similarly, Hoare-Prompt's logic based reasoning improves over basic CoT, however at a smaller rate than TRAILS with accounting for  $\tau = 0.8$ .

#### Finding RQ1

TRAILS improves performance in most of the cases compared to baselines, up to 39% in terms of MCC and 20% for P4. The improvement is particularly striking on the more complex CoCoClaNeL dataset, thanks to the grounding (input, output) pairs that help stirring the reasoning.

## 4.2 RQ2: What is the cost of each step of TRAILS?

In this research question we study the cost of our approach. We give in Table 1 the aggregated cost per task for each approach. Compared to similar approach, TRAILS exhibits a cost overhead in terms of tokens, roughly twice HoarePrompt tokens. We examine the distribution of the tokens cost: in practice, if the phase 2 (output verdict) is on average the most expensive part, it is proportional to the number of inputs generated. Relative to the number of inputs, this part do not contribute much to the overhead. The part that significantly contributes the performance cost is the repair part in phase 1: in practice, especially when the code under test is incorrect with a crashing mistake, the repair can incur half of the total cost. We compared using a MannWhitney test [14] the distribution of tokens for correct and incorrect code samples. For LiveCodeBench, in all case, the p-value was significant at threshold 0.05 with small to medium size effect. As such, correct code assessment would cost a more similar token cost to HoarePrompt's, but the incorrect code assessment led to additional cost. This is nonetheless necessary to ensure the proper validity of the inputs with regard to the task. In the case of CoCoClaNeL, we did not observe statistically significant difference: both correct/incorrect code required repair. Analyzing the log, we note that most of the repair occurs because of mistakes in the input format with the model do not generate the correct amount of input arguments for the task. All CoCoClaNeL tasks requires to inject inputs via standard inputs, with particular constraints. As such, models might not have been widely trained on such input type

and so would need guidance in that step, compared to baselines that do not generate any inputs. Finally, we note that for some models such as Devstral, the LLM struggle to follow the required formatting (inputs, in python block etc.), which further increase the cost. This could, nonetheless, be mitigated via dedicated fine-tuning of the model or optimisation of the given prompt to a particular model using dedicated prompt optimization technique [3].

#### Finding RQ2

TRAILS tokens cost mainly stem from the repair steps incurred both by the necessity to repair inputs even if correct in the case of an incorrect code, as well as particular input format of CoCoClaNeL dataset (standard input). While higher than baselines, TRAILS tokens cost remain of the same order of magnitude than HoarePrompt and could be further optimize via fine-tuning.

## 4.3 RQ3: How stable is the decision of TRAILS across reruns?

In the third research question, we focus on stability of the decision. On top of absolute performance, decision made by an approach should be correctly consistent across runs. That is, rerunning with the same code/specification should not yield a different labels, otherwise it means the approach is sensitive to non-determinism of the model. To verify it, for each approach, we collected tasks across reruns for which: a) all runs labeled the task as "correct" and the code under test was correct, b) all runs labeled the task as "incorrect" and the code under test was incorrect. We then compute the ratio using this number of tasks with all reachable tasks (i.e. task for which one run at least assigned a given label), giving us a score of how "correctly" stable the approach is. A perfect classifier would be consistent across runs while predicting correctly. The naive classifier (always correct/incorrect) would end up with a low score on the naive class given all tasks would be reachable even if it is very stable. Similarly for a random classifier. As such, this score informs on the stability while accounting for possible bias. We give the results are given in Table 2.

From the table, we note that TRAILS is overall the most stable while correctly labeling the code under test, illustrated with the among the highest score across the board compared to the baselines. It is especially more correctly stable when tackling incorrect code with a score of 70%+, meaning it generally consistently assign the

**Table 2: Comparison of verdict stability across two datasets using our approach and baselines. Best result per dataset–metric block is in bold, second best is underlined. 'Correct' is the number of correct code under test for which all runs of one approach labeled it as such (resp 'Incorrect'). We also give the percentage of tasks all labeled 'Correct' (resp. 'Incorrect') which are reachable (i.e. at least one run that labeled the task with this label). For TRAILS, we use  $\tau = 0.8$  which yielded the best performance in RQ1.**

Dataset	Model	TRAILS		HoarePrompt		Zero-shot COT	
		Correct	Incorrect	Correct	Incorrect	Correct	Incorrect
ICB	Qwen3-Coder	57 (60.00%)	53 ( <b>98.15%</b> )	55 (57.89%)	48 ( <u>88.89%</u> )	49 ( <b>63.63%</b> )	59 (81.94%)
	Devstral-Small2	44 (45.36%)	47 ( <b>90.38%</b> )	54 (46.55%)	26 (78.79%)	57 ( <b>49.57%</b> )	29 (85.29%)
	Olmo3.1-Instr	49 ( <b>56.98%</b> )	54 ( <b>85.71%</b> )	43 (45.74%)	43 ( <u>78.18%</u> )	38 ( <u>48.72%</u> )	54 (76.06%)
CoCo	Qwen3-Coder	81 ( <b>35.37%</b> )	59 ( <b>74.68%</b> )	56 (30.43%)	83 ( <u>66.93%</u> )	33 ( <u>32.35%</u> )	115 (55.82%)
	Devstral-Small2	67 ( <u>31.02%</u> )	68 ( <b>73.91%</b> )	70 ( <b>31.39%</b> )	62 ( <u>72.94%</u> )	48 (29.09%)	87 (60.84%)
	Olmo3.1-Instr	100 ( <b>47.17%</b> )	77 ( <b>80.21%</b> )	51 (28.18%)	88 ( <u>69.29%</u> )	37 ( <u>36.63%</u> )	129 (62.32%)

correct label to invalid code and do not assign it to correct code. In comparison, Zero-Shot COT tend to be somewhat correctly stable when dealing with correct code, yet it tends to over assign "incorrect" label to code under test across runs, illustrated by the lower score, 10 to 20 percentage point lower than our approach. It yields similar results when labeling correct code, with a decrease in performance on LiveCodeBench/Olmo (8 percentage point difference). HoarePrompt exhibits a similar patterns: on some models it tends to exhibit a bias pronounced bias towards the "incorrect" labels. For "correct" labels, HoarePrompt is similar to other approaches.

Intuitively, our TRAILS 's grounding via diverse inputs make them less sensitive to variation due to non-determinism. It's less the case for Zero-shot COT or HoarePrompt reasoning which is drastically influenced by the reasoning path taken by the LLM, given there is no grounding inputs, and so more non-determinism [4, 28]. From the user point of view, while Zero-COT or HoarePrompt can yield adequate performance for a lower cost, it comes with at the cost of a lower reliability on a particular individual task. On the contrary TRAILS allows for better stability. The results particularly indicate that this phenomenon is more prevalent on incorrect code, which might lead to detrimental result for the user (determining a code as correct when it's incorrect) in case of an unlucky seed.

**Finding RQ3**

TRAILS shows better correct stability across models/datasets, labeling more consistently the task with the correct label thanks to the grounding inputs. On the contrary, Zero-shot COT and HoarePrompt presents lower stability on tasks due to non-determinism, yielding to a lower decision trust on particular inputs. In particular, the stability is quite lower on incorrect code, leading to a higher possibility of deeming the code correct because of an unlucky seed.

#### 4.4 RQ4: How many consistently and correctly labeled code do each approach label?

Finally, in the last research questions, we study differences in tasks always well labeled (i.e. always correct for correct code and always incorrect for incorrect code) across approaches. The goal is to see

**Table 3: Number of consistently and correctly labeled code by each approach or intersection of approaches across datasets and models. T = TRAILS, H = HoarePrompt, C = Zero-Shot COT**

Model		only T	only H	only C	T∩H	T∩C	H∩C	T∩H∩C	
LiveCodeBench	Qwen3-Coder	Correct	11	2	0	6	2	9	38
		Incorrect	7	2	6	2	9	9	35
	Devstral-Small	Correct	7	3	4	2	4	18	31
		Incorrect	20	3	3	7	10	6	10
	Olmo3.1	Correct	16	4	0	5	4	10	24
		Incorrect	9	2	3	2	12	8	31
CoCoClanL	Qwen3-Coder	Correct	48	12	1	15	3	14	15
		Incorrect	14	5	27	1	11	44	33
	Devstral-small	Correct	25	15	5	16	4	17	22
		Incorrect	22	8	24	7	16	24	23
	Olmo3.1	Correct	57	9	5	14	4	3	25
		Incorrect	8	2	21	2	24	41	43

the number of tasks each approach is better at handling as well as matching behavior between approaches. Results are given in Table 3.

As one can observe, there is a sizable overlap between the approaches ( $T \cap H \cap C$ ) indicating tasks for which, for a given models, all approach consistently assign the correct labels. This behavior is, proportionally, more prominent on LiveCodeBench given the lower number of codes (~70 correct/incorrect codes) than on CoCoClanL (~150 correct/incorrect codes), which is likely because LiveCodeBench's tasks are on average simpler and the code was produced by an LLM. TRAILS generally exposes more unique behavior on LiveCodeBench for both correct/incorrect labels and on CoCoClanL only for correct labels, highlighted by the higher number of tasks that are only labeled consistently by our approach. Such behavior are generally because the code contains long conditions or complex processing. In that case, the model is not burdened by reasoning over the code process itself. We give an example of such task in Listing 1. In that example, both Zero-Shot COT and HoarePrompt will struggles processing the code, hampered by reasoning on the concrete code.

**Listing 1: Correct code consistently misidentified by baselines. (LiveCodeBench - 9ccb)**

```

1 from typing import *
2 class Solution:
3     def removeAlmostEqualCharacters(self, word):
4         result = 0
5         for i in range(len(word)-1):
6             if (i+1)+result >= len(word):
7                 break
8             if abs(ord(word[(i+1)+result])-ord(
9                 word[i+result])) <= 1:
10                result += 1
11        return result
12
13 # HoarePrompt typical output
14 ...
15 { result is the maximum number of consecutive
16   characters in word where each adjacent pair
17   has ASCII values differing by at most 1,
18   starting from some initial position. } #
19   Ignore first stop condition
20 return result
21 { the state is unknown }
22 ...
23
24 # Zero-shot COT typical output
25 ...
26 Let me trace through a simple example:
27 word = "abc"
28 - i=0: check if abs(ord('b')-ord('a')) <= 1 -> abs
29   (98-97) = 1 <= 1, so we increment result to 1
30 - i=1: check if abs(ord('c')-ord('b')) <= 1 ->
31   abs(99-98) = 1 <= 1, so we increment result to
32   2 # Focus on the second condition, when the
33   first one would have been triggered already
34 ...

```

On the other hand of the spectrum, Zero-Shot COT highlights very few unique behavior in terms of correct code, but do so for incorrect code. As seen in previous research question, this might be partially explain by the Zero-shot COT being inclined to label code as incorrect. Zero-Shot COT and HoarePrompt tends to be more similar in consistently labeling correct code (higher number of tasks in the intersections  $H \cap C$ ) in LiveCodeBench and in consistently labeling incorrect code in CoCoClaNeL. In practice, this means there is complementarity in two or three combinations of approaches, meaning they could potentially be combined to further improve the classification performance.

**Finding RQ4**

TRAILS consistently labels more unique code for both correct and incorrect labels, identifying more unique behavior compared to the baselines. That is, TRAILS generally correctly classify code that other approaches struggles with. However,

there are sizeable overlap between all approaches, highlighting complementary between approaches. Zero-Shot COT and HoarePrompt, being pure reasoning, tends to label consistently similar tasks compared to TRAILS.

**5 Approach's limitations**

Despite its advantages, TRAILS can not classify correctly all code under test. We thus dwelled into the potential limitations of the approach. To do so, we manually analyzed code under test for which, on a given dataset, at least two models labeled consistently the code under test with the opposite label using TRAILS. We identified the following main issues:

*Valid inputs causing crash on incorrect code.* We give an example of such case in Listing 2. In that code, there is a bug triggered when  $x > y$  with the increment operation causing a recursion issue. The rest of the code is working correctly ( $x == y$  and  $x < y$ ). Despite this being an identified scenario by TRAILS, given inputs in this scenario will always cause a failure, such inputs are not consider for verification. In that case, the LLM always assert that correct scenarios are correct and the incorrect behavior is never exercised. For such use case, pure reasoning model will generally have an advantage as there are not bounded by concrete execution issues.

**Listing 2: Recursion issue triggered by a non-covered functionality (LiveCodeBench - 2a0f)**

```

1 class Solution:
2     def minimumOperationsToMakeEqual(self, x: int,
3     y: int) -> int:
4         if x <= y:
5             return y - x
6
7         @lru_cache(None)
8         def min_operations(n):
9             if n <= y:
10                return y - n
11            return 1 + min(
12                n % 11 == 0 and min_operations(n
13                // 11) or float('inf'),
14                n % 5 == 0 and min_operations(n //
15                5) or float('inf'),
16                min_operations(n - 1),
17                min_operations(n + 1) # Causes
18                infinite recursion
19            )
20
21        return min_operations(x)

```

*Missing scenario to probe the specification.* We give an example of such case in Listing 3. In that example, the code does not properly account for substring with "s" which starts at another position than 0. Few models run establish such implicit scenario, rather focusing on boundaries for length of "s" or size of "k". Doing so, they occult the bug. When an input from such scenario is generated, models generally identify the error. This signals that, despite prompting for categorical partitioning, this might not be enough to obtain such

**Listing 3: Missing scenario for substring with non 0 starting position (LiveCodeBench - b863)**

```

class Solution:
    def beautifulSubstrings(self, s: str, k: int)
        -> int:
        vowels = set("aeiou")
        n = len(s)
        count = 0
        prefix = defaultdict(int)
        prefix[0] = 1
        v_count = 0

        # Should be a second for loop, from i to n
        for i in range(n):
            if s[i] in vowels:
                v_count += 1
                c_count = (i + 1) - v_count
                if (v_count * c_count) % k == 0 and (
                    v_count - c_count) % 2 == 0:
                    count += prefix[v_count - c_count]
                prefix[v_count] += 1
        return count
    
```

edge case. Looking at reasoning trace of other approaches, inputs examples are close to never used in the reasoning, showing that it is indeed not a trivial scenario to extract, even more so from the specification only.

## 6 Related Works

### 6.1 LLMs for Code, Test and Inputs generations

The use of large language models (LLMs) for code has been a rapidly growing research area in software engineering. Lemieux et al. proposed CodaMOSA, which combines search-based software testing with Codex to escape coverage plateaus by prompting the model for test cases targeting under-covered functions [21]. Rao et al. introduced CAT-LM, a GPT-style model pre-trained with an explicit alignment signal between source code and test files, enabling it to generate contextually coherent tests that outperform much larger general-purpose models [29]. Chen et al. presented ChatUniTest, a framework that uses an adaptive focal context mechanism and a generate-validate-repair loop to produce compilable and coverage-effective unit tests [8]. Dinella et al. propose TOGA [9], a neural approach that reformulates oracle generation as a ranking problem over a grammar-constrained candidate set, handling both exceptional and assertion oracles. Hossain et al. further improved TOGA by introducing TOGLL [16], which fine-tunes a suite of code LLMs with targeted prompts. While showing the capacity of LLMs to generate coherent test and codes, recent studies highlight a misalignment between specification, code and tests [17, 19] which can hamper performance. In particular, findings show that models are prone to asserting the current (potentially buggy) implementation, which limits their fault-detection potential. In contrast, we rely on LLM reasoning using grounded execution via generated test inputs only, limiting the effect of invalid generated oracle. HumanEval+ [23] used LLMs along with manual mutator to create additional test inputs to expand a benchmark. Contrary to us, they rely on

manual mutator which might not preserve complex object nor constraints of the inputs, leading to additional validation. In particular, we showed that repairing the inputs (especially for atypic datasets such as CoCoClaNeL) has the highest impact on cost of the approach. Direct generation with LLMs allow to leverage reasoning to mitigate this effect.

### 6.2 Oracle-Guided Correctness Inference

Complementary to our work with an opposite stance methods, there is a body of literature focusing on oracle-guided correctness inference. Rather than estimating correctness without ground truth, these approaches assume the availability of an existing oracle—typically a reference test suite or a ground-truth solution—and use it to validate, filter, or rank LLM-generated code. The central insight is that passing a sufficiently rich set of pre-existing tests is a reliable proxy for functional correctness, provided the test suite itself is of high quality. AlphaCode [22] operationalises this at scale in competitive programming: it generates millions of candidate programs, filters them against the small set of example test cases provided in the problem statement, and then clusters the surviving candidates before submission. A key limitation of this paradigm is that small test suites systematically overestimate model performance. Mathews et al. [24] shows that supplying public benchmark test cases directly to the LLM during generation—and iterating on failures—improves results on HumanEval and MBPP, with further gains from a repair loop. Similarly, Fakhoury et al. introduce TiCoder [25], a workflow in which a partially formalized intent specification (expressed as user-confirmed test cases) guides iterative code generation. Taken together, oracle-guided approaches differ fundamentally from our method and from oracle-less approaches alike. They require a trusted oracle—either a ground-truth solution, a pre-existing test suite, or a human-confirmed specification—to determine correctness. This makes them highly precise when such an oracle is available, but inapplicable in the typical code evaluation scenario where the oracle is precisely what is missing or in question. Our approach targets this gap: we infer correctness solely from execution traces on generated inputs, without access to ground-truth outputs or any pre-existing test oracle.

### 6.3 Oracle-Less Correctness Inference

The test oracle problem has been recognized as fundamental in software testing. Classic approaches to address it include metamorphic testing, which defines metamorphic relations that relate input/output pairs without requiring explicit oracles, and differential testing, which compares outputs of similar program implementations to detect faults. Yu et al. introduce retromorphic testing, using inverse functions to generate oracles without explicit test cases. Most directly relevant to our work, Valentin et al. propose incoherence, a measure of error estimated without oracles [35]. Fan et al. propose LLMCodeChoice [10], an oracle-guided program selection framework that constructs small distinguishing test suites via differential testing on LLM-sampled programs, using an LLM to predict expected outputs only for those inputs. HoarePrompt [5], which we compare ourselves to, use instead logic based reasoning to infer correctness. Contrary to properties based (metamorphic, intramorphic etc.), we do not require domain knowledge which would require a

proper verification of the property (human or other) contrary to code ran inputs. Approaches such as the one from Valentin et al. requires critically multiple code snippets from the same specification (scale issue) and assume incorrect code will differ across each other which might not hold in all cases. In our case, no mutator is needed and only one sample code is needed, the agreement being based solely on inputs. Finally, compared to HoarePrompt we ground the LLMs reasoning via concrete (input, output) execution, which we show help improve the performance and stability of the method.

## 7 Threats to Validity

**Internal Validity:** Our approach depends on the quality and diversity of generated test inputs. If the categorical partitioning process biases input generation toward certain code patterns, it may not accurately reflect code correctness across diverse specifications, resulting in increased correct code as outlined previously. In practice, however, it offers a better alternative to not generating inputs as we have shown through baselines comparisons. The repair mechanism with a fixed budget may miss valid inputs for complex tasks with unusual requirements or multiple input formats, potentially underestimating the correctness of sophisticated code. This is nonetheless a necessity in order to ensure only valid inputs execute over code. The deduplication step based on code coverage could inadvertently remove important test cases that trigger the same lines but through different logical paths, potentially reducing the effectiveness of outputs verification. This is nonetheless necessary to ensure a minimum number of effective inputs. This issue could be mitigated by adding further information in the deduplication (data-flow process, input size etc.). The binary decision could be affected by LLMs' behavior or non-determinism. In practice, we mitigate this effect by averaging the decision over multiple inputs.

**External Validity:** Our evaluation is limited to two datasets (LiveCodeBench and CoCoClaNeL). LiveCodeBench is a standardized Python benchmark used in multiple studie. CoCoClaNeL is the unique benchmark HoarePrompt tested on, and so we included it for proper comparison. The generalizability to other code generation domains needs to be explored, yet we believe our approach would translate easily. Indeed, the Code requirements step includes extraction of relevant information (mocks for library, resources needed etc.) which would mean the LLM could realistically generates concrete inputs for more advanced tasks. Our evaluation uses three open-source models of mid-range capability that are fairly recent. While results may not generalize to newer model architectures, frontier closed-source models like GPT-4 or Claude, we chose open-source models among the most recent/capable that could run on a single consumer-grade GPU, making the approach deployable in a lot of situations. Finally, we only evaluated on temperature = 0.5 only. Results might vary based on this parameters and certain models could improve/decrease this performance based on it. We chose this value following default values used in baselines.

**Construct Validity:** Our main assumptions for inferring code correctness is that correct code will produce output which the LLMs will recognises as more aligned with the specification on average than for incorrect code. While this assumption is sound, it is impacted by LLM non-determinism. Experiments seem to show the approach is less affected than competing baselines. Moreover,

we offset this effect through averaging over multiple inputs. We use MCC and P4 as primary performance metrics. These metrics are standard metrics used in previous studies when considering binary classifier. The threshold selection ( $\tau = 0.8$ ) was chosen empirically based on our experiments. While it seems to hold well across the board, it may not generalize well to new domains or model variants. In that case, leveraging the calibration mechanism we described might be a good alternative. Non-determinism of runs could affect our results. We mitigated this by rerunning each model/dataset/approach three times.

**Replicability:** We have made sure to describe as thoroughly as possible our steps and processes. All models and datasets used are open-source and available. We also provide a replication package [2].

## 8 Conclusion

This paper addresses the critical challenge of validating large language model-generated code by proposing TRAILS, a novel framework for inferring code correctness through input-grounded reasoning. The fundamental insight of our work is that grounding LLM reasoning with concrete input-output pairs over the specification significantly improves the reliability of correctness assessment. By anchoring reasoning with concrete values rather than allowing models to reason freely over code, we reduce hallucinations and ensure that the specification remains the source of truth throughout the verification process. While TRAILS incurs higher computational costs than baselines due to input generation and repair, these costs remain comparable to formal verification approaches like HoarePrompt and represent a worthwhile trade-off for improved reliability.

Our comprehensive experimental evaluation demonstrates substantial improvements across multiple dimensions. On two benchmarks with three different LLM architectures, TRAILS achieves up to 39% improvement in Matthew Correlation Coefficient and 20% improvement in the P4 metric over Zero-Shot COT reasoning. These gains are made even on CoCoClaNeL, where codes are human written with lower variable semantic, intricate control flow and subtle specifications challenge that pure reasoning approaches might struggle with. Beyond performance, TRAILS exhibits superior correct stability: across multiple runs with different random seeds, our approach consistently produces the same verdict more often than baselines, reducing overall false positives and making the approach more suitable for users.

Future research should explore expanding the approach to more complex scenario, optimizing the input generation and repair process through prompt engineering or model fine-tuning or investigating ensemble combinations of TRAILS with other verification approaches to leverage their complementary strengths. We believe TRAILS represents an important step toward making code generation by LLMs a reliable and trustworthy technology for software development, especially for a growing crowd of non-expert user with limited verification skills.

## Data Availability Statement

Data unavailable during peer-review

References

[1] [n. d.]. Devstral-Small2. <https://devstralsmall2.com/>. Accessed: 2026-03-06.

[2] [n. d.]. ReplicationPackage. Not available during peer-review.

[3] Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, et al. 2025. Gepa: Reflective prompt evolution can outperform reinforcement learning. *arXiv preprint arXiv:2507.19457* (2025).

[4] Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large Language Models for Mathematical Reasoning: Progresses and Challenges. *arXiv:2402.00157* [cs.CL] <https://arxiv.org/abs/2402.00157>

[5] Dimitrios Stamatios Bouras, Yihan Dai, Tairan Wang, Yingfei Xiong, and Sergey Mechtaev. 2025. HoarePrompt: Structural Reasoning About Program Correctness in Natural Language. *arXiv preprint arXiv:2503.19599* (2025).

[6] Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397* (2022).

[7] Ligu Chen, Qi Guo, Hongrui Jia, Zhengran Zeng, Xin Wang, Yijiang Xu, Jian Wu, Yidong Wang, Qing Gao, Jindong Wang, et al. 2024. A survey on evaluating large language models in code generation tasks. *arXiv preprint arXiv:2408.16498* (2024).

[8] Yinghao Chen, Zehao Hu, Chen Zhi, Junxiao Han, Shuiguang Deng, and Jianwei Yin. 2024. Chatunitest: A framework for llm-based test generation. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*. 572–576.

[9] Elizabeth Dinella, Gabriel Ryan, Todd Mytkowicz, and Shuvendu K. Lahiri. 2022. TOGA: A Neural Method for Test Oracle Generation. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. 2130–2141. doi:10.1145/3510003.3510141

[10] Zhiyu Fan, Haifeng Ruan, Sergey Mechtaev, and Abhik Roychoudhury. 2024. Oracle-Guided Program Selection from Large Language Models. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (Vienna, Austria) (ISSTA 2024)*. Association for Computing Machinery, New York, NY, USA, 628–640. doi:10.1145/3650212.3680308

[11] Ahmed Fawzy, Amjed Tahir, and Kelly Blincoe. 2025. Vibe Coding in Practice: Motivations, Challenges, and a Future Outlook—a Grey Literature Review. *arXiv preprint arXiv:2510.00328* (2025).

[12] Molly Q Feldman and Carolyn Jane Anderson. 2024. Non-expert programmers in the generative AI future. In *Proceedings of the 3rd annual meeting of the symposium on human-computer interaction for work*. 1–19.

[13] Francis Geng, Anshul Shah, Haolin Li, Nawab Mulla, Steven Swanson, Gerald Soosai Raj, Daniel Zingaro, and Leo Porter. 2025. Exploring student-AI interactions in vibe coding. *arXiv preprint arXiv:2507.22614* (2025).

[14] Thomas P Hettmansperger and Joseph W McKean. 2010. *Robust nonparametric statistical methods*. CRC press.

[15] Soneya Binta Hossain and Matthew B Dwyer. 2025. Togll: Correct and strong test oracle generation with llms. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE, 1475–1487.

[16] Soneya Binta Hossain and Matthew B. Dwyer. 2025. TOGLL: Correct and Strong Test Oracle Generation with LLMs. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. 1475–1487. doi:10.1109/ICSE55347.2025.00098

[17] Dong Huang, Jie M Zhang, Mark Harman, Mingzhe Du, and Heming Cui. 2024. Measuring the influence of incorrect code on test generation. *arXiv preprint arXiv:2409.09464* (2024).

[18] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. 2024. Qwen2. 5-Coder Technical Report. *arXiv preprint arXiv:2409.12186* (2024).

[19] Michael Konstantinou, Renzo Degiovanni, and Mike Papadakis. 2024. Do llms generate test oracles that capture the actual or the expected program behaviour? *arXiv preprint arXiv:2410.21136* (2024).

[20] Michael Konstantinou, Renzo Degiovanni, Jie M Zhang, Mark Harman, and Mike Papadakis. 2025. YATE: The Role of Test Repair in LLM-Based Unit Test Generation. *arXiv preprint arXiv:2507.18316* (2025).

[21] Caroline Lemieux, Jeevana Priya Inala, Shuvendu K Lahiri, and Siddhartha Sen. 2023. Codamosa: Escaping coverage plateaus in test generation with pre-trained large language models. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 919–931.

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

[23] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in neural information processing systems* 36 (2023), 21558–21572.

[24] Noble Saji Mathews and Meiyappan Nagappan. 2024. Test-Driven Development and LLM-based Code Generation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering (Sacramento, CA, USA) (ASE '24)*. Association for Computing Machinery, New York, NY, USA, 1583–1594. doi:10.1145/3691620.3695527

[25] Noble Saji Mathews and Meiyappan Nagappan. 2024. Test-Driven Development and LLM-based Code Generation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering (Sacramento, CA, USA) (ASE '24)*. Association for Computing Machinery, New York, NY, USA, 1583–1594. doi:10.1145/3691620.3695527

[26] Jaid Naman, Han King, Gu Alex, Li Wen-Ding, Yan Fanjia, Zhang Tianjun, Wang Sida, Solar-Lezama Armando, Sen Koushik, and Stoica Ion. 2024. LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code. *arXiv preprint* (2024).

[27] Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heineman, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, Jacob Morrison, Jake Poznanski, Kyle Lo, Luca Soldaini, Matt Jordan, Mayee Chen, Michael Noukhovitch, Nathan Lambert, Pete Walsh, Pradeep Dasigi, Robert Berry, Saumya Malik, Saurabh Shah, Scott Geng, Shane Arora, Shashank Gupta, Taira Anderson, Teng Xiao, Tyler Murray, Tyler Romero, Victoria Graf, Akari Asai, Akshita Bhagia, Alexander Wettig, Alisa Liu, Aman Rangapur, Chloe Anastasiades, Costa Huang, Dustin Schwenk, Harsh Trivedi, Ian Magnusson, Jaron Lochner, Jiacheng Liu, Lester James V. Miranda, Maarten Sap, Malia Morgan, Michael Schmitz, Michal Guerin, Michael Wilson, Regan Huff, Ronan Le Bras, Rui Xin, Rulin Shao, Sam Skjonsberg, Shannon Zejiang Shen, Shuyue Stella Li, Tucker Wilde, Valentina Pyatkin, Will Merrill, Yapei Chang, Yuling Gu, Zhiyuan Zeng, Ashish Sabharwal, Luke Zettlemoyer, Pang Wei Koh, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. Olmo 3. *arXiv:2512.13961* [cs.CL] <https://arxiv.org/abs/2512.13961>

[28] Nearchos Potamitis, Lars Klein, and Akhil Arora. 2025. ReasonBENCH: Benchmarking the (In) Stability of LLM Reasoning. *arXiv preprint arXiv:2512.07795* (2025).

[29] Nikitha Rao, Kush Jain, Uri Alon, Claire Le Goues, and Vincent J Hellendoorn. 2023. CAT-LM training language models on aligned code and tests. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 409–420.

[30] Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2023. An empirical evaluation of using large language models for automated unit test generation. *IEEE Transactions on Software Engineering* 50, 1 (2023), 85–105.

[31] Lin Shi, Chiyu Ma, Wenhua Liang, Xingjian Diao, Weicheng Ma, and Soroush Vosoughi. 2025. Judging the judges: A systematic study of position bias in llm-as-a-judge. In *Proceedings of the 14th International Joint Conference on Natural Language Processing and the 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*. 292–314.

[32] Mikolaj Sitarz. 2022. Extending F1 metric, probabilistic approach. *arXiv preprint arXiv:2210.11997* (2022).

[33] Philipp Straubinger and Gordon Fraser. 2023. A survey on what developers think about testing. In *2023 IEEE 34th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 80–90.

[34] Qwen Team. 2025. Qwen3 Technical Report. *arXiv:2505.09388* [cs.CL] <https://arxiv.org/abs/2505.09388>

[35] Thomas Valentin, Ardi Madadi, Gaetano Sapia, and Marcel Böhme. 2025. Incoherence as Oracle-less Measure of Error in LLM-Based Code Generation. *arXiv preprint arXiv:2507.00057* (2025).

[36] Zhilong Wang, Lan Zhang, Chen Cao, Nanqing Luo, Xinzhi Luo, and Peng Liu. 2023. How Does Naming Affect LLMs on Code Analysis Tasks? *arXiv preprint arXiv:2307.12488* (2023).

[37] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[38] Lin Yang, Chen Yang, Shutao Gao, Weijing Wang, Bo Wang, Qihao Zhu, Xiao Chu, Jianyi Zhou, Guangtai Liang, Qianxiang Wang, et al. 2024. On the evaluation of large language models in unit test generation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*. 1607–1619.

[39] G Udny Yule. 1912. On the methods of measuring association between two attributes. *Journal of the Royal Statistical Society* 75, 6 (1912), 579–652.