RETHINKING LOGIC IN AI: A NOVEL BENCHMARK INSPIRED BY POLYNOMIAL ANALOGUE OF GANDY'S FIXED POINT THEOREM

Anonymous authors

Paper under double-blind review

Abstract

This paper introduces a novel benchmark for evaluating the logical reasoning capabilities of Large Language Models (LLMs), grounded in the polynomial analogue of Gandy's classical fixed point theorem. Since this theorem can be used to describe the P-complete HornSAT problem, and our benchmark is based on this theorem, our benchmark thus covers all problems from class P and shows that serious problems have already arisen in this class, not to mention those benchmarks whose complexity classes are NP-complete and NP-hard. Drawing on concepts from mathematical logic, we design a parameterized set of recursively definable problems where the objective is for LLMs to predict whether a problem belongs to an inductively definable set of polynomial complexity. By varying the parameters, we generate problem instances of differing complexity. Our experiments reveal that current state-of-the-art LLMs with zero-shots promts fail to reliably solve even the most straightforward cases despite an effective deterministic algorithm existing. Even advanced models like GPT-4 exhibit significant biases in solving benchmark problems. These findings highlight the limitations of modern LLMs as code interpreters, even in basic scenarios, and underscore the necessity for hybrid LLM/interpreter systems. Furthermore, they emphasize the importance of developing quantitative tests for reasoning, given the increasing reliance on LLM-based systems in decision-making applications.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

1 INTRODUCTION

034 1.1 MOTIVATION

The question of whether large language models (LLMs) can think has been a subject of intense debate for many years (Mirzadeht et al., 2024), and it is only getting hotter as large language models 037 solve more and more problems. To address this question, researchers are designing increasingly complex tests and challenges for LLMs, which allow us to assess the cognitive capabilities of these models more precisely. From a logical perspective, a key component of any intelligent system being 040 able to think is the ability to reason and construct proofs from basic axioms, as well as the ability 041 to work with recursive definitions (Goncharov & Nechesov, 2021b) If we look at the problems from 042 the other side, all problems can be divided into complexity classes P, NP-complete, NP-hard, etc. 043 Many benchmarks try to cover all of these classes at once(Lizhou Fan et al., 2024), but if LLMs 044 cannot solve problems from class P well, then it is pointless to try to study questions in the more complex classes NP-complete and NP-hard. These findings are supported by the results obtained from most studies that aim to tackle problems spanning across different complexity levels simulta-046 neously. When we approach problems from the perspective of formalization, the majority of them 047 can be articulated through non-recursive descriptions with first-order logic formulas. However, when 048 a problem permits a recursive representation, its formulation in logical terms becomes significantly more compact. It is convenient to work with such recursive descriptions of problems and send them to the input of LLMs together with the objects for which these problems need to be solved. 051

In this paper, we propose to focus on a unique combination of problems that lie in the class P and at the same time admit a recursive description. If LLMs cannot solve these problems, then most often it is pointless to move on to more complex problems from the NP-complete and NP-hard

054 classes, since their complexity becomes higher and at the same time most of them also admit a 055 recursive description. To navigate through recursive structures while remaining within the confines 056 of class P, the PAG theorem (Appendix A) and FPAG theorem (Nechesov & Goncharov, 2024) prove 057 invaluable, enabling us to define sets and functions that can be described recursively with polynomial 058 complexity. In this investigation, we focus on the PAG theorem, demonstrating that even within this framework, LLMs encounter significant difficulties in determining whether an object belongs to an inductively defined collection of entities. The uniqueness of the benchmark is achieved due to its 060 unique properties - in essence, we are trying to solve recursively defined problems, but at the same 061 time do not jump out of the class P. 062

063 The HornSAT problem (Adebayo et al., 2022), which lies in the class P and is a simplification of 064 the NP-complete SAT problem, also falls under the recursive description. The HornSAT problem can be solved using a modified PAG theorem, where we can to transfer the main concepts of the 065 PAG theorem from nested lists to the set of horn formulas, and also had to abandon generating 066 families in favor of the set of horn formulas. This was achieved by redefining the function γ , which 067 now takes some horn formula as input and produces another shorter horn formula according to a 068 recursive P-computable algorithm for Horn satisfiability (hor). The operator Γ^{HF} acts on the set 069 of all subsets of horn formulas in the set of all subsets of horn formulas. In this case, the smallest fixed point Γ^* of the operator Γ^{HF} is P-computable. Without loss of generality, we can assume 071 that the HornSAT problem is a special case of a modified polynomial analogue of Gandy's theorem. 072 A similar approach can be used to describe the 2-SAT problem using a modified PAG theorem. 073 Since the 2-SAT problem reduces to a $2N \times 2N$ adjacency matrix, a transitive closure is iteratively 074 constructed for it using a modified version of the Floyd-Warshall algorithm. (2sa). Due to the fact 075 that the HornSAT problem is p-complete, we can assume that our benchmark based on the PAG theorem covers all problems from the class P. This is its uniqueness and universality. 076

Large language models have a transformative effect on various domains of science, including language processing and understanding (Ouyang et al., 2022), music and speech processing (Agostinelli et al., 2023), smart cities and digital twins of complex objects, and even drug and protein design (Madani et al., 2023). These models, such as GPT-4 (Achiam et al., 2023), Claude (cla) and their open-source competitors like Llama (Dubey et al., 2024) and Mistral (Jiang et al., 2023), have demonstrated remarkable capabilities in tasks ranging from simple text generation to complex mathematical problem-solving. The remarkable advancements of LLMs were enabled by the abundance of unlabelled text data (Penedo et al., 2024) and the effective unsupervised training at scale (Kaplan et al., 2020).

Despite these successes, there are notable limitations in the current methodologies for evaluating the
 reasoning capabilities of LLMs. Standardized benchmarks (Hendrycks et al., 2020; Yue et al., 2024;
 Zheng et al., 2023) have shown that LLMs excel in few-shot and zero-shot tasks, often outperforming
 previous state-of-the-art models by large margins. However, these benchmarks may not fully capture
 the models' ability to generalize and reason logically, as evidenced by various studies highlighting
 often unexpected failures Nezhurina et al. (2024) that contradict the claimed strong capabilities.

Addressing these limitations is crucial for applying LLMs in real-world scenarios, where accurate,
 logical reasoning is essential. For instance, in fields like healthcare, legal services, and automated
 decision-making, the reliability of AI models can have significant implications for outcomes and
 trustworthiness. Therefore, there is an urgent need for new evaluation frameworks Gao et al. (2023);
 Guha et al. (2024); Wang (2024) that can more accurately assess the logical capabilities of LLMs in
 different settings.

098 In response to this need, we introduce a new benchmark designed to evaluate the reasoning abilities 099 of LLMs in recursive problems. Our approach involves generating [object, condition] pairs, where objects are nested lists of elements (for example, numbers) and conditions are expressed with recur-100 sive functions. The task of the LLM is to determine whether the object satisfies all conditions. This 101 setup, expressed in Python language, leverages the extension of fixed-point Gandy's theorem (Gon-102 charov & Nechesov, 2021b), ensuring that the check can be performed efficiently at most quadratic 103 times relative to the object's size. Despite an efficient verification algorithm, our findings indicate 104 that all tested LLMs fail to solve even simple instances of these problems, highlighting a signif-105 icant gap in their logical reasoning capabilities. The source code of the benchmark is available 106 at https://anonymous.4open.science/r/logic-lm-0B9C.

108 2 RELATED WORK

As the capabilities of large language models themselves increase, so does the number of various benchmarks that test their capabilities in various areas including the ability to solve mathematical and logical problems. Most often, the problems that arise when testing LLMs can be divided into complexity classes: P, NP-complete, and NP-hard.

In our research, we primarily focus on problems from the P class, as we rely on the PAG (Gon-charov & Nechesov, 2021b) and FPAG theorems (Nechesov & Goncharov, 2024) the powerful tools for constructing polynomially computable sets of inductively defined objects and also for defining recursive functions of polynomial complexity. This allows us to see how the larger language models handle these constructions.

The complexity arises from LLMs' need to effectively handle recursive operations that inevitably emerge when tackling such problems. Recursive definitions serve as a key factor in increasing complexity. By carefully managing recursion, we can either maintain operations within the P class or transition to NP-complete or NP-hard classes.

As our research shows, even within the complexity class P, there exist tasks that pose a challenge for large language models, suggesting that higher-level classes such as NP-complete and NP-hard become simply unsolvable in most cases. It is futile to engage in experimentation with issues in these domains unless one possesses a comprehensive understanding of the root causes for tasks from class P.

In the work under consideration Lizhou Fan et al. (2024), attention is directed towards tasks spanning various complexity levels, encompassing P, NP-complete, and NP-hard classes. The research aims to comprehensively explore these diverse complexity classes, rather than focusing narrowly on a specific aspect. However as we noted earlier in the context of this work, where researchers try to solve problems belonging to the NP-complete and NP-hard classes without having fully comprehended the intricacies of problems within the class P.

In another work (Rishi Hazra et al., 2024) the solution of the 3-SAT problem using LLM occurred with a small number of free variables (no more than 10). That is a 3-SAT problem in which it is possible to find a solution manually by enumeration. It turned out that even with such a small number of parameters, LLMs cope poorly with them. As expected, with an increase in parameters, their capabilities sharply decrease.

If we look at benchmarks that have made a lot of noise in the LLM community, then of course this
is the Alice problem. (Nezhurina et al., 2024) This problem can be solved by an 8-year-old child,
but most of the leading LLMs could not do it. If we want to build trusted AI algorithms based on
LLMs, then it is very important to start with the simplest problems to understand where they make
mistakes and where they do not. The most acceptable complexity class for these problems is in class
P. That is why it is so important to study this complexity class in as much detail as possible, which
is what we do in this work.

147 148

3 BACKGROUND

149

159

160

The idea of our benchmark is to use a parametric class of problems that can be efficiently generated and solved, yet the complexity of individual problem instances can be tuned. We revert to the Polynomial Analogue of Gandy's Fixed Point Theorem (Goncharov & Nechesov, 2021b), or PAG-theorem, to define such a class. In this section, we describe how to define this class of problems in simple terms and refer the readers interested in stricter descriptions to Appendix A.

Let us introduce notation for a class of polynomially solvable problems \mathcal{F} . Each problem consists of a pair (P, x), where $P \in \mathcal{F}$ is a recursively defined boolean function (a predicate), and x is a candidate object. The objective is to determine whether or not x belongs to P, e.g. if

- $P(x) == 1 \tag{1}$
- here 1 in Eq. 1 means logical True, and 0 means False, the symbol == means a comparison operation.

The objects $x \in \mathcal{X}$ are arbitrarily nested lists of elements e. To be concrete, in this work, we consider lists of natural numbers $e \in \mathcal{N}$, but any elements can be used. We inductively denote a nested list xas follows:

$$x = \langle x_0, \dots, x_N \rangle$$

$$x_k = e_k \quad \text{or} \quad \langle x_{k,0}, \dots, x_{k,M} \rangle, \quad k \in [0, \dots, N]$$
(2)

165 166 167

168

173 174

175 176 We use composite indices of length d to describe elements x at the d-th level of nesting.

The boolean functions in \mathcal{F} must adhere to certain structures imposed by the PAG theorem. Consider a set of predicates $\{P_0(x), \ldots, P_I(x)\}$. As objects x are lists, the structure of $P \in \mathcal{F}$ will depend on how it is applied to individual elements of the list. In this work we consider boolean functions on lists $x = \langle x_0, \ldots, x_N \rangle$ of the following form:

 $\Phi_{j}(x_{0}, \dots, x_{N}) = Q_{j,0}(x_{0}) * \dots * Q_{j,N}(x_{N})$ $Q_{j,i}(x_{i}) = \begin{cases} P_{j,i}(x_{i}) \\ (x_{i} == e_{i}) \end{cases}$ (3)

The symbol * denotes a generic logical operator, which can be either a conjunction (and, \wedge) or a disjunction (or, \vee). The function $\Phi_k(x)$ is thus a logical formula, where predicates $P_{j,i} \in \mathcal{F}$ are applied to individual elements of x.

As nested lists $x_k \in \mathcal{X}$ may have different number of elements at each level of nesting, we need to specify how the action of $P \in \mathcal{F}$ depends on the length of $x_k = \langle x_{k,0}, \ldots, x_{k,N_k} \rangle$. We do so by introducing the operators $\gamma_i[x], i \in [0, \ldots, I]$:

$$\gamma_i[x_k] = \begin{cases} \Phi_{N_k}(x_{k,0}, \dots, x_{k,N_k}), & \text{if } \gamma_i(x_k) \text{ is defined for length } N_k \\ 0, & \text{otherwise} \end{cases}$$
(4)

Based on the length of x_k , the operators γ_i either produce some logical formula Φ_{N_k} (which in turn consists of some possibly recursive logical functions $P_{k,l} \in \mathcal{F}$) or False. We denote the set of all predicates that can be produced by γ_i as F_i , which is also called a generating family. Note that $F_i(x_k)$ is also (a possibly recursive) logical function. Given $x_i \in \mathcal{X}$

190 191

192

199

203 204

205

206 207

208

183 184 185

$$F_i(x_k) = \gamma_i[x_k] = \begin{cases} \Phi_{N_k}(x_{k,0}, \dots, x_{k,N_k}) = Q_{N_k,0}(x_0) * \dots * Q_{N_k,N_k}(x_{N_k}) \\ 0 \end{cases}$$
(5)

The generating family $F_i(x_k)$ in Eq. 5 is a recursive function if $Q_{N_k,l}$ may take a value of $F_i(x_{k,l})$, for example.

We are now ready to formulate the PAG theorem. Given a set of polynomially computable nonrecursive boolean functions $\mathcal{P} = \{P_0, \dots, P_J\}$, given a set of boolean functions $\mathcal{F} = \{F_0, \dots, F_I\}$ with the structure defined by Eq. 5, where

$$Q_{N_k,l}(x) = F_{k,l}(x) \in \mathcal{F} \text{ or } P_{k,l}(x) \in \mathcal{P} \text{ or } (x_l == e_l)$$
(6)

and each $\Phi_{N_k}(x_0, \dots, x_{N_k})$ may contain no more than one copy of $F_k(x_j), \forall k \in [0, \dots, I]$ for each $x_j, j \in [0, \dots, N_k]$ in its definition. If $x \in \mathcal{X}$ is a list with finite depth and width, then the equation

$$F_i(x) == 1 \tag{7}$$

can be decided in a polynomial number of steps for any $F_i \in \mathcal{F}$. We will use this result to build a test system for large language models.

4 Methods

Based on the results of the PAG theorem, we build a generator of problem instances with varying complexity. Each problem has the form of a set of generating families \mathcal{F} , a test object $x \in \mathcal{X}$ and the equality $F_i(x) == 1, F_i \in \mathcal{F}$. We call the set of generating functions a *condition* and the test object x the *probe*. In the following sections, we will show how to express the *condition* as a Python program and present a probe generator algorithm, which allows the building of probes of any depth and any desired result when substituted in Eq. 7. The resulting Python code can be directly supplied to LLMs, and the response can be evaluated without human intervention, thus providing an automatic benchmark system.

216 4.1 CONDITION GENERATOR

218 The setup of the PAG theorem has to be mapped to a Python program to make a condition generator. 219 We start by expressing logical formulas Φ in Eq. 3. Consider an input list with N = 3 elements x = $\langle x_0, x_1, x_2 \rangle$ of size 3. We first need to specify element-wise functions $Q_i, i \in [0, 1, 2]$. According 220 to the requirements of the PAG-theorem, they can be either constant comparisons, non-recursive predicates \mathcal{P} or possibly recursive generating families \mathcal{F} (see Eq. 6). For simplicity, we do not use 222 non-recursive predicates \mathcal{P} here (an example of such predicate would be a Python function **bool** (), which returns True for any non-zero input or some other non-recursive function). We define I = 2224 generating families called is_member_0() and is_member_1(). An example of the function Φ_N 225 for N = 3 is shown below: 226

 $res = is_member_0(x[0])$ and (x[1] == 2) and $is_member_1(x[2])$

227 228

229 230

231

232

249 250

265

266 267

268

269

To define a generating family, one must map inputs x to different Φ functions based on the input length. An example of a set of two generative families called is_member_0() and is_member_1 () is shown below. Notice that these functions may be mutually recursive.

Listing 1: Φ_3 function example

```
233
     1 def is_member_0(x):
234
     2
           if len(x) == 2:
235
               res = (x[0] != 15) and (x[1] != 61)
     3
     4
           elif len(x) == 3:
236
               res = (is_member_1(x[0])) and (is_member_0(x[1])) and (x[2] ==
237
           49)
238
           else:
     6
239
               return False
     7
240
           return res
     8
241
     1 def is_member_1(x):
242
     2
           if len(x) == 2:
243
               res = (is_member_1(x[0])) and (is_member_0(x[1]))
     3
244
           elif len(x) == 3:
     4
245
               res = (x[0] != 46) and (x[1] == 95) and (x[2] != 7)
     5
246
     6
           else:
247
     7
               return False
           return res
     8
248
```



251 We made several design choices to translate the requirements of the PAG theorem into concrete Python functions. Our code generator produces a set of I recursive functions as shown in Lst. 2. Each function has if-blocks, which map inputs starting from length $N_i = 2$ to length N_i^{max} to a 253 corresponding logical formula Q_{N_i} . We use only and operators in logical formulas. Every logical 254 formula can contain up to $b_i, b_i \leq N_i$ generating families $F_i \in \mathcal{F}$ in its definition, which is selected 255 at random; the rest of the entries in Q_{N_i} are defined as argument comparisons to random integer 256 constants. Finally, each generating family F_i has a logical formula of length N_i^T , which does not 257 contain any calls to generating formulas $F_j \in \mathcal{F}$, e.g., completely composed of comparisons of 258 the arguments to constants. In the example in Lst. 2 $N_0^T = 2$ for $F_0 = is_member_0$ () and 259 accordingly $N_1^T = 3$ for $F_1 = \text{is_member_1}$ (). We found that having a terminal block in every 260 $F_i \in \mathcal{F}$ is necessary for solving solutions of different nesting depths in Eq. 7. 261

We implemented the condition generator as described above using the template library Jinja2 (Ronacher, 2008). The generator produces random Python functions acting on lists x. The second half of the benchmark system is a probe generator, which is described in the next section.

4.2 PROBE GENERATOR

In this section, we devise an algorithm to build nested lists $x \in X$ satisfying Eq. 7. Notice that there are usually multiple solutions. For example, both $\times 0$ and $\times 1$ in Lst. 3

x0 = [-39, 21]

	Listing 3: Solution examples	
	Lisung 5. Solution examples	
return numb need	True when passed to is_member_0() from Lst. 2, but the evaluation will take a difference of steps due to the different depth of the lists. To have an efficient benchmark, one work to generate $x \in \mathcal{X}$ with a depth of nesting d , which will yield a desired result in Eq. 7. Such a efficiently built	ent uld h x
To de	e concluding bunk. Escribe the algorithm, let us first consider a logical formula $Q_{N_i}^T$ completely made for compari- guments to constants. We call this formula a <i>terminal</i> formula.	ing
	$Q_{N_i}^T = (x_0 == e_0) \wedge \dots \wedge (x_{N_i} == e_{N_i})$	(8)
It is e	easy to find the argument which satisfies Q_{T}^{T} :	. ,
10 15 0	$x = \langle o_1, \dots, o_m \rangle$	(0)
By cł be co	hanging a single entry in Eq. 9, one can build a non-solution. In general, <i>terminal</i> formulas of mposed of non-recursive predicates $P_j \in \mathcal{P}$ provided it is possible to solve $P_j(x) = 1$.	can
If a g consi	generating family F_i produces $Q_{N_i}^T$, then x from Eq. 9 is also a solution to $F_i(x) = 1$. Note a generating family F_j , which produces a logical formula Q_{N_k} containing F_i :	ow
	$Q_{N_k} = F_i(x_m) \land \widetilde{Q}_{N_k - 1}(x_0, \dots, x_{m-1}, x_{m+1}, \dots, x_{N_k}) \tag{1}$	10)
here	\widetilde{Q}_{N_k-1} is a subformula of Q_{N_k} without the term $F_i(x_m)$. If \widetilde{x}	=
$\langle \tilde{e}_0, .$	$\ldots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \ldots, \tilde{e}_{N_k}$ satisfies \tilde{Q}_{N_k-1} , then	
$\langle \tilde{e}_0, .$	$\dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \widetilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_k} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle $	11)
$\langle \tilde{e}_0, .$ is a sean alg	$ \begin{array}{l} \ldots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \ldots, \tilde{e}_{m-1}, \langle e_0, \ldots, e_{N_i} \rangle, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \end{array} $ olution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1.	11) 'ith
$\langle \tilde{e}_0, .$ is a sean alg	$ \begin{array}{l} \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \\ \text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ \text{gorithm outlined in Alg. 1.} \\ \hline \textbf{rithm 1} \text{ Probe generator} \\ \hline \end{array} $	11) 'ith
$\langle \tilde{e}_0, .$ is a sum of an algorithm of the second seco	$ \begin{array}{l} \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, e_{\tilde{N}_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, e_{\tilde{N}_k} \rangle \qquad ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed we gorithm outlined in Alg. 1.} \\ \hline \mathbf{rithm 1 Probe generator} \\ \hline \mathbf{unction SOLVETERMINAL}(Q^T, r) \qquad \triangleright \text{ Solves } Q^T(x) = x \\ \hline \mathbf{r} \langle e_{N_k}, e_{N_k} \rangle = x \\ \hline \mathbf{r} \langle e_{N_k} \rangle = x \\ \hline \mathbf{r} \langle e_{N_k}, $	$\frac{11}{rith} = r$
$\langle \tilde{e}_0, .$ is a sean algorithm of $Algorithm for the second s$	$ \begin{array}{l} \ldots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \ldots, \tilde{e}_{m-1}, \langle e_0, \ldots, e_{N_i} \rangle, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \qquad ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ \text{gorithm outlined in Alg. 1.} \end{array} $	$\frac{11}{rith} = r$
$\langle \tilde{e}_0,$ is a sea an alg $\frac{\text{Algori}}{1: \text{ fr}}$ 2: r 3: e	$ \begin{array}{ll} \ldots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0, \ldots, \tilde{e}_{m-1}, \langle e_0, \ldots, e_{N_i} \rangle, \tilde{e}_{m+1}, \ldots, e_{\tilde{N}_k} \rangle \qquad ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ \text{gorithm outlined in Alg. 1.} \end{array} $	$\frac{11}{r}$
$\langle \tilde{e}_0, .$ is a sean algorithm of A lgorithm of A lgorithmode A lgorithm of A lgori	$ \begin{array}{l} \ldots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \ldots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y = \langle \tilde{e}_0, \ldots, \tilde{e}_{m-1}, \langle e_0, \ldots, e_{N_i} \rangle, \tilde{e}_{m+1}, \ldots, \tilde{e}_{N_k} \rangle \qquad ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ \text{gorithm outlined in Alg. 1.} \\ \hline \textbf{rithm 1 Probe generator} \\ \hline \textbf{unction SolveTerminal}(Q^T, r) \\ x \leftarrow \langle e_0, \ldots, e_N \rangle, \text{ s.t. } Q^T(x) = r \\ \textbf{eturn } x \\ \textbf{nd function} \end{array} $	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a sean alg $\overline{\text{Algorithmatrix}}$ 1: ff 2: r 3: e 4: f	$\begin{split} \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x$ $x \leftarrow \langle e_0, \dots, e_N \rangle, \text{ s.t. } Q^T(x) = r$ eturn x nd function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = x$	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a sean alg Algon 1: ff 2: 3: e 4: ff 5:	$\begin{split}, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x \leftarrow \langle e_0, \dots, e_N \rangle$, s.t. $Q^T(x) = r$ eturn x rid function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = $ if $d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula Q	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a sea an alg $\overline{\text{Algor}}$ 1: ff 2: r 3: e 4: ff 5: 6:	$\begin{split}, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x \leftarrow \langle e_0, \dots, e_N \rangle$, s.t. $Q^T(x) = r$ eturn x rid function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = $ if $d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula $Q^T \in F_i$	$\frac{11}{r}$
$\langle \tilde{e}_0, .$ is a so an alg $\overline{\text{Algor}}$ 1: ff 2: 7: 3: e 4: ff 5: 6: 7: 2	$\begin{split}, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x$ $x \leftarrow \langle e_0, \dots, e_N \rangle$, s.t. $Q^T(x) = r$ eturn x end function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = $ if $d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula $Q^T \in F_i$ $x \leftarrow SOLVETERMINAL(Q^T, r)$	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a sean alg Algor 1: ft 2: r 3: e 4: ft 5: 6: 7: 8: 0:	$\begin{split}, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x$ $x \leftarrow \langle e_0, \dots, e_N \rangle, \text{ s.t. } Q^T(x) = r$ eturn x eturn x function GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = $ if $d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula $Q^T \in F_i$ $x \leftarrow \text{SOLVETERMINAL}(Q^T, r)$ return x return x	$\frac{11}{r}$
$\langle \tilde{e}_0, .$ is a second state of a second st	$\begin{split}, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \end{split} $ (colution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x$ $x \leftarrow \langle e_0, \dots, e_N \rangle, \text{ s.t. } Q^T(x) = r$ eturn x and function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = $ if $d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula $Q^T \in F_i$ $x \leftarrow$ SOLVETERMINAL (Q^T, r) return x else $\triangleright x$ is a nested list, recurse into a random non-terminal formula Q^T Eind a random $Q \in F$ at $Q = F_0 \land Q^T \land Q^T$	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a solution and solution of the second state of	$\begin{split} &, \tilde{e}_{m-1}, \tilde{e}_{m+1},, e_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0,, \tilde{e}_{m-1}, \langle e_0,, e_{N_i} \rangle, \tilde{e}_{m+1},, e_{N_k} \rangle ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ & \text{gorithm outlined in Alg. 1.} \\ \hline \\ $	11) rith = r = r Q^T Q_k
$\langle \tilde{e}_0,$ is a sea an alg Algor 1: f 2: r 3: e 4: f 5: 6: 7: 8: 9: 10: 11:	$\begin{split} \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ y &= \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed we gorithm outlined in Alg. 1.} \\ \hline \text{rithm 1 Probe generator} \\ \hline \text{unction SOLVETERMINAL}(Q^T, r) & \triangleright \text{ Solves } Q^T(x) = x \\ & \forall \in \langle e_0, \dots, e_N \rangle, \text{ s.t. } Q^T(x) = r \\ \hline \text{eturn } x \\ \text{end function} \\ \hline \text{unction GENERATESOLUTION}(\mathcal{F}, F_i \in \mathcal{F}, d, r) & \triangleright \text{ Finds } x \text{ of nesting level } d, \text{ s.t. } F_i(x) = x \\ \hline \text{if } d = 0 \text{ then} & \triangleright x \text{ is a shallow list, solve terminal formula } 0 \\ \hline \text{Find } Q^T \in F_i \\ x \leftarrow \text{SOLVETERMINAL}(Q^T, r) \\ \hline \text{return } x \\ \hline \text{else} & \triangleright x \text{ is a nested list, recurse into a random non-terminal formula } 0 \\ \hline \text{Find a random } Q_k \in F_i, \text{ s.t. } Q_k = F_{j_0} \wedge \dots \wedge F_{j_b} \wedge \tilde{Q}_k \\ y \leftarrow \text{SOLVETERMINAL}(\tilde{Q}_k, 1) & \triangleright \text{ Solve a non-recursive part of the form } 0 \\ \hline \text{For a solution } D \\ \hline \text{Solve a non-recursive part of the form } D \\ \hline Solve a non-recursive part of the $	11) vith = r = r Q^T Q_k ula
$\langle \tilde{e}_0,$ is a sea an alg Algon 1: f 2: r 3: e 4: f 5: 6: 7: 8: 9: 10: 11: 12:	$\begin{array}{llllllllllllllllllllllllllllllllllll$	$\frac{11}{r}$
$\langle \tilde{e}_0,$ is a sea an alg Algon 1: f 2: r 3: e 4: f 5: 6: 7: 8: 9: 10: 11: 12: 13:	$\begin{split} &, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, e_{\tilde{N}_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, e_{\tilde{N}_k} \rangle & ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ & \text{gorithm outlined in Alg. 1.} \\ \hline \text{rithm 1 Probe generator} \\ \hline \text{unction SOLVETERMINAL}(Q^T, r) & \triangleright \text{ Solves } Q^T(x) = \\ & x \leftarrow \langle e_0, \dots, e_N \rangle, \text{ s.t. } Q^T(x) = r \\ & \text{eturn } x \\ & \text{nd function} \\ \hline \text{unction GENERATESOLUTION}(\mathcal{F}, F_i \in \mathcal{F}, d, r) & \triangleright \text{ Finds } x \text{ of nesting level } d, \text{ s.t. } F_i(x) = \\ & \text{if } d = 0 \text{ then} \\ & \vdash x \text{ is a shallow list, solve terminal formula } 0 \\ & \text{Find } Q^T \in F_i \\ & x \leftarrow \text{SOLVETERMINAL}(Q^T, r) \\ & \text{return } x \\ & \text{else} \\ & \vdash x \text{ is a nested list, recurse into a random non-terminal formula } 0 \\ & \text{Find a random } Q_k \in F_i, \text{ s.t. } Q_k = F_{j_0} \wedge \cdots \wedge F_{j_k} \wedge \tilde{Q}_k \\ & y \leftarrow \text{SOLVETERMINAL}(\tilde{Q}_k, 1) \\ & \vdash \text{Solve a non-recursive part of the form } 0 \\ & \text{For } j \in [j_0, \dots, j_b] \text{ do } \\ & y_j \leftarrow \text{GENERATESOLUTION}(\mathcal{F}, F_j, d-1, r) \\ & \text{ord form } \\ \end{array}$	11) rith rith = r Q^T Q_k ula
$\langle \tilde{e}_0,$ is a solution of a set	$\begin{split} & \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, e\tilde{N}_k \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, e\tilde{N}_k \rangle \end{split} $ ((olution of depth one of $F_j(x) = 1$. In general, the solution of depth d can be constructed we gorithm outlined in Alg. 1. rithm 1 Probe generator unction SOLVETERMINAL (Q^T, r) \triangleright Solves $Q^T(x) = x$ $x \leftarrow \langle e_0, \dots, e_N \rangle$, s.t. $Q^T(x) = r$ eturn x nd function unction GENERATESOLUTION $(\mathcal{F}, F_i \in \mathcal{F}, d, r)$ \triangleright Finds x of nesting level d , s.t. $F_i(x) = if d = 0$ then $\triangleright x$ is a shallow list, solve terminal formula G Find $Q^T \in F_i$ $x \leftarrow \text{SOLVETERMINAL}(Q^T, r)$ return x else $\triangleright x$ is a nested list, recurse into a random non-terminal formula G Find a random $Q_k \in F_i$, s.t. $Q_k = F_{j_0} \land \dots \land F_{j_b} \land \tilde{Q}_k$ $y \leftarrow \text{SOLVETERMINAL}(\tilde{Q}_k, 1)$ \triangleright Solve a non-recursive part of the form for $j \in [j_0, \dots, j_b]$ do \triangleright Breadth-first search-like iteration $y_j \leftarrow \text{GENERATESOLUTION}(\mathcal{F}, F_j, d - 1, r)$ end for $\pi \cdot (x)$	11) with $rac{r}{r}$ = r Q_{r} Q_{k} ula
$\langle \tilde{e}_0,$ is a sea an alg Algor 1: f 2: r 3: e 4: f 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16:	$\begin{split} & \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle & ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ & \text{gorithm outlined in Alg. 1.} \end{split}$	11) with rith = r Q^T Q_k ula ion
$\langle \tilde{e}_0,$ is a sea an alg Algor 1: f 2: r 3: e 4: f 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17:	$\begin{split} & \dots, \tilde{e}_{m-1}, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \text{ satisfies } \tilde{Q}_{N_k-1}, \text{ then} \\ & y = \langle \tilde{e}_0, \dots, \tilde{e}_{m-1}, \langle e_0, \dots, e_{N_i} \rangle, \tilde{e}_{m+1}, \dots, \tilde{e}_{N_k} \rangle \qquad ((\text{olution of depth one of } F_j(x) = 1. \text{ In general, the solution of depth } d \text{ can be constructed w} \\ & \text{gorithm outlined in Alg. 1.} \end{split}$	11) rith = r = r Q^{T} Q_{k} ula

The algorithm 1 is a breadth-first search-like recursive algorithm. Recall that every $F_i \in \mathcal{F}$ has a terminal formula Q_i^T by our design choice a; hence we can always generate a shallow solution $x = \langle e_0, \ldots, e_N \rangle$ such that $F_i(x) = 1$. If we have to generate a nested solution of depth d, then we choose a non-terminal formula $Q_k \in F_i, Q_k = F_{j_0} \land \cdots \land F_{j_b} \land \tilde{Q}_k$ at random, solve a nonrecursive part \tilde{Q}_k of Q_k , and recurse into corresponding generating families $F_j, j \in [j_0, \ldots, j_b]$ until a required depth d is reached. Notice that when generating non-solutions $\bar{x}, F_i(\bar{x}) = 0$, we choose not to satisfy only the terminal formulas at depth d. The latter guarantees that in order to verify $F_i(\bar{x}) == 1$, the evaluator will take at least d steps to check the nested elements of \bar{x} .

324 4.3 EVALUATION 325

Having introduced the algorithms for data generation, we need to specify how they will be supplied
 to language models. In this work, we considered two kinds of models: conversational LMs and the
 models trained specifically for code completion. For conversational models, we used the following
 prompt:

```
330
              What is the result of the following Python code?
              '''python
331
              {condition}
332
              x = \{probe\}
333
              print(is_member_0(x))
334
              ...
335
              Answer only 'True' or 'False'
336
                                     Listing 4: Prompt for conversational LMs
337
338
        In Lst. 4 the placeholder condition takes the definitions of the generative families \mathcal{F}, and probe takes
339
        the value of x. For code-specific models, there is no clear way to ask the model to evaluate some
340
        expressions. Instead, we relied on the ability of these models to enforce the correctness of Python
341
        expressions:
342
      1 {condition}
343
      2 \mathbf{x} = \{\text{probe}\}
344
      3 assert is_member_0(x) ==
345
                                    Listing 5: Prompt for code completion LMs
346
347
        The response of the language models is easy to verify automatically using the proposed prompts.
348
349
         5
             RESULTS
350
351
352
         5.1 Setup
353
         We selected several state-of-the-art general purpose and code completion LMs for the experimental
354
        evaluation listed in Tab. 1. In the experiments, we compared the accuracy of LMs concerning each
355
356
                                  Table 1: Tested models along with their versions
357
358
                        Conversational LMs
                                                                           Code completion LMs
359
            GPT-3.5 [turbo-0125] (Ouyang et al., 2022)
                                                                     PolyCoder [2.7B] (Xu et al., 2022)
360
          GPT-4-turbo [2024-04-09] (Achiam et al., 2023)
                                                             CodeLlama [7B-Python-hf] (Roziere et al., 2023)
361
                     GPT-40 [2024-05-13] (gpt)
                                                                  Starcoder-2 [3B] (Lozhkov et al., 2024)
                 Llama-3 [8B] (Dubey et al., 2024)
                                                                Stable Code [3B] (Pinnaparaju et al., 2024)
362
364
         parameter of the condition and probe generators. The combinations of parameters we used are listed
        in Tab. 2.
365
366
        Table 2: Parameters of the example generator used in experiments. I - total number of functions,
367
         N_i^{max} - number of if-blocks in each function, b_i - maximal number of calls to recursive functions in
368
        each if-block (branching number), N_i^T - size of the terminal if-block, d - nesting depth of the test
369
        list x.
370
371
                                      Nimax
                                                     Λ
                                                           d
                                                                Varying parameters
                                                b_i
372
                                         \frac{1}{2}
                                                      2
                                                           2
                                                                    I, N_i^{max}, d
                                 2
                                                1
373
                                 11
                                         2
                                                      2
                                                           2
                                                1
                                                                        b_i
                                        12
                                                      2
                                                           2
                                                                        N_i^T
374
                                 2
                                                1
375
```

376

We generated 160 random problem instances with $F_i(x) = 1$ (*positive* instances) and the same number of instances with $F_i(x) = 0$ (*negative* instances) for each parameter combination of the generator. We found that with this sample size, the standard deviation of all metrics is within 10%.
Each model was tested using the same data. Additionally, we used the Bernoulli distribution as a baseline.

We used four binary classification metrics to compare the behavior of language models in the experiments: true positive rate (sensitivity, TPR = true positives/total predicted positives), true negative rate (specificity, TNR = true negatives/total predicted negatives), balanced accuracy BA = $\frac{\text{TPR}+\text{TNR}}{2}$ and the Youden's J statistic J = TPR + TNR - 1. Youden's statistic estimates the probability of an informed decision.

387 388

389

399

419 420

```
5.2 MODEL COMPARISON
```

We observed poor performance of every tested LLM in the benchmark across a wide range of gener-390 ator parameters. Even for the simplest problem instances Lst. 6, balanced accuracy never exceeded 391 75%, despite the relative simplicity of the underlying verification process. Note that as each answer 392 is either True or False, a simple random guess already achieves 50% accuracy. The accuracy of 393 most small-scale models rapidly approached 50% with the increase of the probe's depth, the number 394 of functions, and the number of if-blocks, with low sensitivity to the length of the terminal block. We 395 also found that the accuracy of stronger models (GPT-3.5 and Codellama) increases with the number 396 of recursive branches in functions $F_i(x)$, presumable because it makes easier to verify the probe x 397 by not checking the result of most of the if-blocks. The accuracy's dependence on the benchmark 398 parameters is provided in Fig. 1.

```
What is the result of the following Python code?
400
            '''python
401
     2
            {def is_member_0(x):
     3
402
     4
                if len(x) == 2:
403
                    res = (x[0] != -20) and (x[1] != 1)
     5
404
                elif len(x) == 3:
     6
                    res = (is_member_1(x[0])) and (x[1] == 20) and (x[2] == 30)
405
     7
     8
                else:
406
                     return False
     9
407
     10
                return res
408
           def is_member_1(x):
     11
409
     12
                if len(x) == 2:
                    res = (x[0] == 65) and (is_member_0(x[1]))
410
     13
                elif len(x) == 3:
     14
411
                    res = (x[0] != 24) and (x[1] != -15) and (x[2] == 37)
     15
412
     16
                else:
413
                     return False
     17
414
                return res}
     18
           x = \{ [3, -44] \}
415
     19
           print(is_member_0(x))
     20
416
     21
417
           Answer only 'True' or 'False'
418
                             Listing 6: An example set of a simple problem
```

421 Surprisingly, we found that an advanced model GPT-4-turbo performed worse than random and 422 worse than its simpler versions, see Fig. 1. The observation that GPT-4-turbo is strongly biased on 423 benchmark problems motivated us to check the behavior of LMs on positive and negative instances 424 separately (TPR and TNR metrics). We found that every model in the benchmark demonstrated 425 a significantly biased behavior on "True" and "False" subsets of problems, with some models 426 (PolyCoder, StarCoder) generating a constant answer regardless of the condition used. GPT-3.5 427 demonstrated the slightest bias of all tested LLMs, followed by LLama3 and GPT-40. Again, GPT-428 4-turbo demonstrates a strange bias towards incorrect answers. An example of model behavior with 429 varying probe depth on separate problem subsets is shown in Fig. 2, and additional examples can be found in Appendix B. In all experiments, Youden's J statistic did not exceed 0.45, with typical 430 values of around 0.2, suggesting that LLMs can barely use the condition part of the problem to 431 predict the answer.





Figure 2: TPR, TNR and Youden's *J* statistic dependence on the depth of probe *x*. Despite having almost perfect TPR, StarCoder, Polycoder, and StableCode have zero TNR, suggesting a highly biased behavior.

540 6 DISCUSSION

The construction of benchmark systems is essential steps in the evaluation of LLMs performance.
Our reasoning benchmark is based on the results of the polynomial analog of Gandy's classical fixed
point theorem and provides means for principled LLMs evaluation. We found that even smallest
problem instances are challenging for state of the art LLMs.

The unexpected finding was a significant bias inherent to every tested LM. No model demonstrated an expected Bernoulli distribution-like behavior when tested on uniform positive or negative problem instances, with GPT-3.5 coming closer to the expected distribution. This observation suggests that LLMs may have difficulty with the generation of random data. As the most extreme case, biased behavior renders GPT-4-turbo the worst performer in our test, suggesting that bias increases with the model scale.

Contrary to what may be expected, models trained to understand code demonstrated worse performance than general-purpose LLMs. We plan to adapt our benchmark to natural language formulation to make it even more natural for general purpose LLMs.

Overall, our experiments demonstrate an unusual behavior of LLMs for recursive codeunderstanding tasks. None of the models can consistently solve even simple problem instances despite the existence of a deterministic algorithm with quadratic complexity (Depth-First search with backtracking). Our findings stipulate the revision of language model training procedures to eliminate biases we observed, particularly in contexts where high accuracy and reliability in reasoning are required.

REFERENCES

562

563

569

578

579

580

- 564 565 566 566 564 Floyd-warshall algorithm. https://en.wikipedia.org/wiki/Floyd-Warshall_ algorithm. Accessed: 2024-11-12.
- The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/
 Model_Card_Claude_3.pdf. Accessed: 2024-09-30.
- GPT-40 contributors. https://openai.com/gpt-40-contributions. Accessed: 2024-09-30.
- 572 573 HornSAT algorithm. https://en.wikipedia.org/wiki/Horn-satisfiability. Accessed: 2024-11-12.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
 - Adebayo, Sathasivam, Ali, and MKM. Hornsat solver using agent-based modelling in hopfield network. *Intelligent Systems Modeling and Simulation II. Studies in Systems, Decision and Control, vol 444. Springer, Cham,* 2022.
- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon,
 Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating
 music from text. *arXiv preprint arXiv:2301.11325*, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The Ilama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Yurii Leonidovich Ershov. *Definability and computability*. Springer Science & Business Media, 1996.
- Yanjun Gao, Dmitriy Dligach, Timothy Miller, John Caskey, Brihat Sharma, Matthew M Churpek,
 and Majid Afshar. Dr. bench: Diagnostic reasoning benchmark for clinical natural language
 processing. *Journal of biomedical informatics*, 138:104286, 2023.

601

621

622

623

624

- 594 Sergey Goncharov and Andrey Nechesov. Solution of the problem p=1. *Mathematics*, 10(1):113, 595 2021a. 596
- Sergey Goncharov and Andrey Nechesov. Polynomial analogue of gandy's fixed point theorem. 597 Mathematics, 9(17):2102, 2021b. 598
- Sergey Goncharov and Andrey Nechesov. Semantic programming for ai and robotics. In 2022 IEEE 600 International Multi-Conference on Engineering, Computer and Information Sciences (SIBIR-CON), pp. 810-815. IEEE, 2022. 602
- 603 Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, et al. Legalbench: A collabo-604 ratively built benchmark for measuring legal reasoning in large language models. Advances in 605 Neural Information Processing Systems, 36, 2024. 606
- 607 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and 608 Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint 609 arXiv:2009.03300, 2020. 610
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 611 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 612 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 613
- 614 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, 615 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language 616 models. arXiv preprint arXiv:2001.08361, 2020. 617
- Wenyue Hua Lizhou Fan, Haoyang Ling Lingyao Li, and Yongfeng Zhang. Nphardeval: Dynamic 618 benchmark on reasoning ability of large language models via complexity classes. arXiv preprint 619 arXiv:2312.14890, 2024. 620
 - Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The next generation. arXiv preprint arXiv:2402.19173, 2024.
- Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton, 625 Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language models 626 generate functional protein sequences across diverse families. Nature Biotechnology, 41(8):1099-627 1106, 2023. 628
- 629 Iman Mirzadeht, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad 630 Farajtabart. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large 631 language models. arXiv preprint arXiv:2410.05229, 2024.
- Andrey Nechesov. Semantic programming and polynomially computable representations. Siberian 633 Advances in Mathematics, 33(1):66-85, 2023. 634
- 635 Andrey Nechesov and Sergey Goncharov. Functional variant of polynomial analogue of gandy's 636 fixed point theorem. Mathematics 2024, 12(21), 3429, 2024. 637
- 638 Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. Alice in wonderland: Simple tasks showing complete reasoning breakdown in state-of-the-art large language models. 639 arXiv preprint arXiv:2406.02061, 2024. 640
- 641 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 642 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-643 low instructions with human feedback. Advances in neural information processing systems, 35: 644 27730-27744, 2022. 645
- Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro 646 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data 647 at scale. arXiv preprint arXiv:2406.17557, 2024.

- 648 Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, Ashish Datta, 649 Maksym Zhuravinskyi, Dakota Mahan, Marco Bellagente, Carlos Riquelme, et al. Stable code 650 technical report. arXiv preprint arXiv:2404.01226, 2024. 651 Gabriele Venturato Rishi Hazra, Pedro Zuidberg Dos Martires, and Luc De Raedt. Can large lan-652 guage models reason? a characterization via 3-sat. arXiv preprint arXiv:2408.07215, 2024. 653 654 Armin Ronacher. Jinja2 documentation. Welcome to Jinja2—Jinja2 Documentation (2.8-dev), 2008. 655 656 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for 657 code. arXiv preprint arXiv:2308.12950, 2023. 658 659 Zeyu Wang. Causalbench: A comprehensive benchmark for evaluating causal reasoning capabilities 660 of large language models. In Proceedings of the 10th SIGHAN Workshop on Chinese Language 661 Processing (SIGHAN-10), pp. 143–151, 2024. 662 Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of 663 large language models of code. In Proceedings of the 6th ACM SIGPLAN International Sympo-664 sium on Machine Programming, pp. 1–10, 2022. 665 666 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, 667 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multi-668 modal understanding and reasoning benchmark for expert agi. In Proceedings of the IEEE/CVF 669 Conference on Computer Vision and Pattern Recognition, pp. 9556–9567, 2024. 670
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
 chatbot arena. Advances in Neural Information Processing Systems, 36:46595–46623, 2023.
 - A THE POLYNOMIAL ANALOGUE OF GANDY'S FIXED POINT THEOREM

The polynomial analog of Gandy's fixed point theorem (PAG-theorem) (Goncharov & Nechesov, 2021b; Nechesov, 2023) is based on the classical Gandy's fixed point theorem (Ershov, 1996), which is important not only in mathematical logic but also in programming. The classical Gandy's theorem states that if we have a fixed model:

$$\Omega = \langle \mathcal{N}, 0, +, \cdot, \leq \rangle$$
 of the signature $\sigma_0 = \langle 0, +^2, \cdot^2, \leq^2 \rangle$

⁶⁸³ Denote by σ^* the extension of the signature σ_0 obtained by adding symbols for all Σ -functions (Ershov, 1996) on Ω and constant symbols for all elements from \mathcal{N} , and let Ω^* denote the corresponding enrichment of Ω .

Define the operator $\Gamma_{\Phi[\overline{x}]}^{\Omega^*}$ as next:

674 675

676

681 682

686

687 688 689

694

695

697 698

$$\Gamma^{\Omega^*}_{\Phi[\overline{x}]}(Q) = \{(e_1, \dots, e_{k-1}) \mid \langle \Omega^*, Q \rangle \models \Phi(e_0, \dots, e_{k-1})\} \cup Q$$
(12)

where \models symbol means checking the truth of the formula $\Phi(e_0, \ldots, e_{k-1})$ with initialized variables on the model $\langle \Omega^*, Q \rangle$, and $Q \subseteq \mathcal{N}^k$ is a truth set of an extendable predicate, $(e_1, \ldots, e_{k-1}) \in Q$.

693 We associate the next sequence of the sets of the truth of predicates:

 $\Gamma_0, \Gamma_1, \ldots, \Gamma_\alpha$ where α is an ordinal.

696 with the monotone operator $\Gamma_{\Phi[\overline{x}]}^{\Omega^*}$ as follows:

$$\Gamma_0 = \emptyset, \Gamma_{\alpha+1} = \Gamma_{\Phi[\overline{x}]}^{\Omega^*}(\Gamma_{\alpha})$$
 for none limit ordinal, ..., $\Gamma_{\alpha} = \bigcup_{\beta < \alpha} \Gamma_{\beta}$ for limit ordinal

Let α be the smallest ordinal such that $\Gamma_{\alpha+1} = \Gamma_{\alpha}$ then $\Gamma^* = \Gamma_{\alpha}$ is a smallest fixed point.

700 **Theorem 1** (Gandy's fixed point theorem)

Let $\Phi(P^+)$ be a Σ -formula of the signature $\sigma^* \cup \{P^{(k)}\}$ in which the predicate symbol P enters

702 703 704	positively and x_0, \ldots, x_{k-1} be a list of the different free variables of the formula Φ then the smallest fixed point Γ^* of the operator $\Gamma_{\Phi[\overline{x}]}^{\Omega^*}$ is a Σ -predicate on Ω^* .
705 706 707	In the PAG-theorem, several conditions from the classical Gandy's fixed point theorem were mod- ified or strengthened so that the fixed point of the operator remains a polynomially computable (P-computable) predicate.
708 709 710 711	To prove this, a P-computable hereditary-finite list superstructure $HW(\mathfrak{M})$ of signature σ was cho- sen as the basic model (Goncharov & Nechesov, 2021b). The base set $HW(M)$ contains hereditary- finite lists, which are inductively constructed from the basic elements of the base set M of the model \mathfrak{M} of the signature $\sigma_0 \subset \sigma$. The basic list operations for $HW(\mathfrak{M})$ are:
712	• $head(x)$ - operation of taking the last element of a list x
714	• $tail(x)$ - operation of removing the last element from a list x
715	• $cons(x, y)$ - operation of adding to the end of a list x of a list y
716 717	• $cons(x, y)$ operation of adding to the end of a first x of a first y
718	• $conc(x, y)$ - concatentation operation of two fists x and y respectively
719	There are also relations:
720	• $x \in u_{-}$ "x is an element of u"
722	$x \in \mathcal{G}^{-1}$ is an element of \mathcal{G}^{-1}
723	• $x \in y$ - "x is an initial segment of y"
724 725	Formulas of the first-order logic are considered only with the bounded quantifiers $\forall x \in t, \exists x \in t, \forall x \subseteq t, \exists x \subseteq t, where t is a standard term.$
726 727 728	In order to apply the PAG theorem generally, it is necessary to construct a P-computable GNF system based on the following components:
729	• A finite alphabet Σ .
730 731	• An extended alphabet Ω .
732	• A special logical language L where L-formulas and L-programs are defined
733	 P computable hereditary finite list superstructure HW(M) of the signature σ
734	Ficher and the second s
736	• Finite sets of extendable predicates P_1, \ldots, P_n .
737	• Generating families of formulas F_{P_1}, \ldots, F_{P_n} .
738 739	• P-computable functions $\gamma_1, \ldots, \gamma_n$ that, given an element of HW(M), construct suitable generating formulas or return False.
740 741	• and other conditions from (Goncharov & Nechesov, 2021b)
742	Theorem 2 (Polynomial analogue of Gandy's fixed point theorem)
743	Let G be a p-computable GNF-system then the smallest fixed point Γ^* of the operator $\Gamma^{HW(\mathfrak{M})}_{-}$
744	is P-computable.
745	The results obtained in the PAG theorem allow us to inductively define sets of objects of varying
747	complexity using generating families of L-formulas and, at the same time, guarantee their recog-
748	nizability in polynomial time. Inductivity upwards essentially defines recursivity downwards, i.e.,
749	the algorithm that checks whether an object belongs to a class of objects, working with various list
750	sistency with one of the families of generating formulas. This entire process can be programmed so
751 752	that the question of whether an element belongs to a set is solved in polynomial time (Goncharov & Nechesov, 2021a: 2022)
the second second second	1 (OCHOSOV, 2021u, 2022).

Corollary 1 Let G is a P-computable GNF-system. Let the computational complexity of all basic functions, predicates, and functions $\gamma_i(x)$ in G be at most $O(|x|^p)$. Then the computational complexity of the smallest fixed point Γ^* is at most $O(|x|^{p+2})$.

B ADDITIONAL DATA



Here, we provide additional benchmark data. Pairs of TPR and TNR plots for various benchmark parameters are listed in Fig. 3.

Figure 3: Pairs of TPR and TNR plots concerning different parameters of the benchmark Youden's statistic behavior depending on the parameters of the benchmark is shown in Fig. 4

