# Assessing the Reverse-Engineering Abilities of Large Language Models

#### Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031 032 033

034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Using AI to create autonomous researchers has the potential to accelerate scientific discovery. A prerequisite for this vision is analyzing whether an AI model can identify the underlying structure of a black-box system from its behavior. In this paper, we explore how well a large language model (LLM) learns to identify a black-box function from passively observed versus actively collected data. We investigate the reverse-engineering capabilities of LLMs across three distinct types of black-box systems, each chosen to represent different problem domains where future autonomous AI researchers may have considerable impact: programs, formal languages, and math equations. Through extensive experiments, we show that LLMs fail to extract information from observations, reaching a performance plateau that falls short of the Bayesian inference ideal. However, we demonstrate that prompting LLMs to not only observe but also intervene—actively querying the black-box with specific inputs to observe the resulting output—improves performance by allowing LLMs to test edge cases and refine their beliefs. By providing the intervention data from one LLM to another, we show that this improvement is partially tethered to the process of generating effective interventions, paralleling results in the literature on human learning. Further analysis reveals that engaging in interventions can help LLMs escape from two common failure modes: overcomplication, where the LLM falsely assumes prior knowledge about the blackbox, and *overlooking*, where the LLM fails to incorporate observations. These insights provide practical guidance for helping LLMs more effectively reverseengineer black-box systems, supporting their use in making new discoveries.

#### 1 Introduction

Developing intelligent systems to accelerate scientific discovery has been a long-standing goal of artificial intelligence research (Gil et al., 2014; Wang et al., 2023). Despite rapid progress in creating large language models (LLMs) for understanding text and solving problems in math and coding, automating science poses a different kind of challenge. A core aspect of scientific discovery is reverse-engineering the mechanism behind a black-box system, which requires capabilities beyond responding to a one-off query. In particular, reverse-engineering often involves 1) understanding observed data in order to develop hypotheses, 2) designing experiments to actively acquire informative data from the black-box to test those hypotheses, and 3) describing and communicating the results.

Existing work using LLMs for automating scientific processes either focuses on static observational data (Rmus et al., 2025; Shojaee et al., 2025) or emulates scientific workflows using "LLM scientists" with many moving parts (Gandhi et al., 2025; Schmidgall et al., 2025). In contrast, research in related fields has used carefully controlled tasks to evaluate whether machine learning systems can perform key aspects of reverse-engineering, including inductive reasoning (Rule et al., 2024), learning causal features from passive data (Lampinen et al., 2023), and optimal experimental design (Chaloner & Verdinelli, 1995; Rainforth et al., 2024). This work is often informed by work in cognitive science, which has studied how humans engage in active learning using methods in which the source (either passive observation or active experimentation) and content of data can be differentiated (Markant & Gureckis, 2010; 2014). Such controlled methods have not been applied to state-of-the-art LLMs, leaving fundamental questions unanswered: "How well can LLMs make inferences from passive observations?" and "Can they actively collect data to refine their hypothe-ses?".

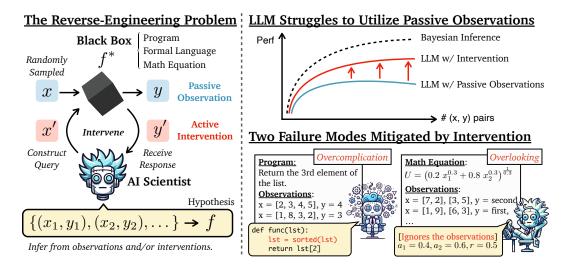


Figure 1: **Reverse-engineering.** Left: Defining the problem. The AI scientist will obtain passive observations from the black box or collect data through active intervention to construct a hypothesis. Right (top): with only passive observations, the LLM cannot make effective use of the data and lags behind Bayesian inference; allowing the LLM to intervene improves performance. Right (bottom): effective intervention can mitigate two common failure modes: overcomplication and overlooking.

To answer these questions, we systematically study LLMs on three reverse-engineering tasks inspired by the cognitive-science literature and selected to mimic challenges in scientific settings: reconstructing list-mapping programs (Rule et al., 2024), formal languages (McCoy & Griffiths, 2023), and math equations (Foster et al., 2019). Through extensive experiments, we show that LLMs are limited in their ability to make inferences from observations, leading to performance plateaus compared to Bayesian models. However, allowing LLMs to perform interventions—generating test cases or queries to collect new, informative data—can significantly improve their performance.

Through further experiments where outcomes of interventions conducted by one LLM become observational data for another, we show that the benefits of intervention seem to come from the LLM testing and refining its own beliefs rather than simply collecting higher-quality data. This is similar to a phenomenon observed in human learning, where people show limited benefit from interventions generated by others (Markant & Gureckis, 2010; 2014). Further investigation reveals that generating interventions seems to help LLMs overcome two failure modes: 1) *overcomplication*, where the LLM tends to construct overly-complex hypotheses, and 2) *overlooking*, where the LLM neglects observations or draws overly-generic conclusions without careful checking.

#### Our contributions are as follows:

- Drawing inspiration from controlled studies of human cognition, we formalize *reverse-engineering* as a core problem for assessing the scientific discovery capabilities of LLMs and design three black-box tasks to facilitate such an assessment.
- We demonstrate empirically that frontier LLMs still struggle, relative to Bayesian inference, at reverse-engineering these black boxes from only passive observations.
- We show that LLMs can perform interventions to obtain more informative data, and that effective intervention mitigates the failure modes of *overcomplication* and *overlooking*.
- We show that performance degrades when repurposing the LLM's intervention data as observations, pinpointing the mechanism behind the improvements it produces and highlighting a potential pitfall for exchanging knowledge among LLMs.

## 2 RELATED WORK

**Inductive Inference** Some of the earliest work on reverse-engineering appears under the label of *inductive inference* for "hypothesizing a general rule from examples" (Angluin & Smith, 1983).

Classic instances of this problem include work on identifying the underlying structure of a finitestate automaton through observations of its input-output behavior (Rivest & Schapire, 1987; 1989). While this problem typically considers passive observations, seminal work on active learning focuses on analyzing the benefits of actively querying inputs to solicit the most-informative outputs from the unknown function of interest (Littlestone, 1988; Angluin, 1988; Settles, 2009). The key distinction between these seminal works and ours is the attention towards LLMs and assessing their capacity for successfully identifying different types of black boxes from input-output examples.

Bayesian Optimal Experiment Design An adjacent line of work considers the sequential de-

sign of experiments which maximally yield information gain about an unknown parameter of interest (Lindley, 1956; DeGroot, 1962; Chaloner & Verdinelli, 1995; Rainforth et al., 2024); one may interpret these methods as studying a non-LLM-focused, Bayesian analogue of the reverse-engineering problem we formulate in the subsequent section, where a learner begins with a prior distribution over the black box in question and must maximally reduce epistemic uncertainty (Der Kiureghian & Ditlevsen, 2009) with a given budget of experiments. To the extent that LLMs may implicitly engage with an underlying approximate posterior inference scheme (Xie et al., 2021; Griffiths et al., 2024; Zhu & Griffiths, 2024a; Falck et al., 2024; McCoy et al., 2024), the reverse-engineering capabilities studied in this work can be tied to this Bayesian optimal experiment design problem.

125 126 127

128

129

130

131

132

133

134

135

136

137

138

139

140

141 142

143

144

145

108

109

110

111

112

113

114 115 116

117

118

119

120

121

122

123

124

**Reinforcement Learning** The fundamentals of the reverse-engineering problem also connect with various ideas studied in the context of reinforcement learning (RL) (Sutton & Barto, 1998). Any model-based RL agent (Sutton, 1990; 1991; Brafman & Tennenholtz, 2002; Strehl & Littman, 2008) naturally engages with a particular instance of the reverse-engineering problem where the black-box function in question is the transition function and/or reward function of a Markov Decision Process (MDP) (Bellman, 1957; Puterman, 1994). The distinction explored in this work between a LLM that passively observes versus actively intervenes on the black box in question has a direct connection to the exploration challenge in RL, which has profound impact on an agent's ability to recover an accurate model of the world (Thrun & Möller, 1991; Deisenroth & Rasmussen, 2011; Strens, 2000; Osband et al., 2013); while recent work (Arumugam & Griffiths, 2025) has studied how to improve exploration with LLMs, this paper focuses on assessing the innate capabilities of LLMs to actively query informative data. Ostrovski et al. (2021) demonstrate the ineffectiveness of passive learning with deep RL agents and their need to intervene so as to correct misunderstandings about the world; our work provides the LLM complement to their findings. The KWIK learning framework of Li et al. (2008) provides a theoretical analysis for reverse-engineering a MDP transition function when a learner must either confidently estimate the environment dynamics or say "I don't know" (Walsh et al., 2009; Li & Littman, 2010; Sayedi et al., 2010; Szita & Szepesvári, 2011; Abernethy et al., 2013). Finally, there is a connection between intervention for effective reverseengineering and meta RL (Liu et al., 2021), with recent work showing that passive learning can be effective with LLMs once there is an effective exploration strategy capable of yielding high-quality observations (Lampinen et al., 2023); naturally, the latter problem is precisely what we demonstrate interventions allow LLMs to solve for themselves in reverse-engineering tasks.

146 147 148

149

150

151

152

153

154

155

156

157

158

**LLMs for Automating the Scientific Process** With the rapid advances in LLMs, recent work has explored using them to automate different parts of the scientific process such as ideation (Si et al., 2024), assistance (Gottweis et al., 2025), writing research papers (Lu et al., 2024; Starace et al., 2025), or emulating AI scientists in simulated environments (Schmidgall et al., 2025). Additionally, multi-modal and multi-agent AI models have driven significant progress in applications such as protein science (O'Neill et al., 2025), while frameworks like MatPolit (Ni et al., 2024) integrate human cognitive insights to accelerate discoveries in materials science. These works utilize the abundant knowledge stored in LLMs to directly tackle real-world complexity in science (Reddy & Shojaee, 2025). Recent work also shows that LLMs can autonomously generate and test hypotheses to advance automated scientific discovery (Agarwal et al., 2025). However, the complexity of these settings and the resulting agents make it hard to disentangle the consequences of the engineering choices that go into these systems. Our work focuses on simple and controllable black boxes to study the core capabilities of the LLMs themselves.

159 160 161

**Understanding Failure Modes in LLMs** Recently, many works have examined the failure modes (Aggarwal et al., 2025) of formal reasoning in LLMs. It has been observed that LLMs can exhibit failure modes of both "overthinking" (Chen et al., 2024) and "underthinking" (Wang et al., 2025) when tackling mathematical problems and code generation (He et al., 2025; Sprague et al., 2024) and (Cuadron et al., 2025; Sprague et al., 2024; Sui et al., 2025; Cemri et al., 2025). To understand LLM abilities beyond formal reasoning tasks, recent work has leveraged insights and datasets from cognitive science (Frank, 2023; Binz & Schulz, 2023; Coda-Forno et al., 2024; Ying et al., 2025). In particular, researchers have started to use cognitive science to explore the failed behaviors in LLMs (Ku et al., 2025). Using these methods, researchers have found that LLMs sometimes overestimate human rationality (Liu et al., 2024a), exhibit inconsistencies in probability judgments (Zhu & Griffiths, 2024b), and perform worse as a result of engaging in reasoning (Liu et al., 2024b). In a similar vein, our work draws upon cognitive science to design the black boxes used in our reverse-engineering experiments.

## 3 REVERSE ENGINEERING

#### 3.1 PROBLEM FORMULATION

We define a black box  $f^*: \mathcal{X} \to \mathcal{Y}$  as a deterministic function that maps a query  $x \in \mathcal{X}$  to a response  $y \in \mathcal{Y}$  through its internal dynamics. The **reverse-engineering** problem is for a model to infer the internals of a black box  $f^*$  (such as list mapping programs, production rules of formal languages, and math equations) from a sequence of query-response pairs  $\mathcal{O} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \subset$  $\mathcal{X} \times \mathcal{Y}$  (Figure 1). We consider two cases of the reverse-engineering problem: **observation-only** and **observation-intervention**. In the observation-only scenario, all the queries are randomly sampled from  $\mathcal{X}$  and the corresponding response  $y_i = f^*(x_i)$  is generated by the black box from a uniform distribution to construct the observation set. A large language model  ${\cal M}$  must generate a hypothesis  $f = \mathcal{M}(\mathcal{O})$  without further interaction with the black box. This setting assesses the model's ability to perform inductive reasoning (Angluin & Smith, 1983). In the observation-intervention scenario, the LLM is first given a set of observations  $\mathcal{O}$  obtained in the observation-only scenario and is instructed to interact with the black box in a multi-round fashion. In each round, the LLM chooses one of the following actions: 1) construct a new query  $x_{N+1}$  to query the black box and obtain the response  $y_{N+1}$ , 2) construct a new query-response pair  $(x_{N+1}, y'_{N+1})$  and check its validity using the black box  $(\mathbb{1}[y'_{N+1} = f^*(x_{N+1})])$ , or 3) stop and conclude with a hypothesis f about the black box. Before constructing the new query, the LLM can analyze the current oservations with strategies such as verbalizing its current belief or describing the current hypothesis (§5.2). Before the LLM chooses to stop or reaches the maximal number of rounds, the query-response pairs obtained during intervention are appended to  $\mathcal{O}$  for the next round.

## 3.2 Black-Box Types

Drawing on the literature on inductive inference in cognitive science, we select tasks commonly used to study learning of complex relationships to design our black-box systems and scale them up for evaluation with LLMs. These three distinct black-box function classes – Program, Formal Language, and Mathematical Equation – simulate problems encountered in scientific reverse-engineering scenarios. Due to space constraints, detailed black-box designs are relegated to Appendix B.

**Program.** We use list-mapping programs (Rule et al., 2024) for the Program black-box. Each program implements a lambda expression (e.g., (lambda (singleton (third \$0)))) in Python, where the query is a list of integers and the response is an integer.

**Formal Language.** The Formal Language black-box is defined by a simple program that generates sequences of symbols. For example, the language  $A^nB^n$  generates sequences consisting of some number of As followed by the same number of Bs. The black-box allows the LLM to intervene by validating if a string is allowable under the rule. We define 46 distinct black boxes each based on a language from Yang & Piantadosi (2022) or McCoy & Griffiths (2023).

**Math Equation.** We use the Constant Elasticity of Substitution (CES) formulation from economics (Foster et al., 2019) as the Math Equation black-box. The utility  $U = (\sum_i a_i x_i^r)^{\frac{1}{r}}$  is given by the weights  $a_i$ , the ratio r, and the quantities of each kind of goods  $x_i$ . The LLM queries the black box with two lists of item types with quantities. The response says which list has higher utility.

#### 3.3 EVALUATION PROTOCOL

A black-box can be represented in multiple ways, rendering evaluation challenging. For example, two black-boxes can be compared through their descriptions in natural language (descriptive evaluation) or whether they respond similarly to the same queries (functional evaluation; see §J). In this paper we focus on **descriptive evaluation**, where the black-box  $f_{\rm NL}^*$  is expressed in natural language, due to its communicative nature and real-world use (Chopra et al., 2019; Gandhi et al., 2025). The LLM-generated hypothesis  $f_{\rm NL}$  is scored by an LLM judge against the black-box on a 0-10 scale based on the criteria of each black-box type ( ${\tt score} = {\tt LM-Judge}(f_{\rm NL}, f_{\rm NL}^*)$ ). We use descriptive evaluation for Program and Formal Language. As the Math Equation does not require verbalization beyond the weights and ratio, we report the flipped root mean square error (1 - RMSE) between the inferred parameters and ground truth.

## 4 EXPERIMENTS

**Experimental setup.** We use different versions of GPT-4o (Hurst et al., 2024) for reverse-engineering (gpt-4o-2024-08-06, dubbed as reverse-engineer LLM) and as the judge (gpt-4o-2024-05-13, dubbed as the judge LLM). We use greedy decoding of both the reverse-engineer and the judge LLMs and report performance over 3 seeds. For the observation-only experiments, we report performance for number of observations  $N = \{2, 10, 15, 20, 30, 60\}$ . For the observation-intervention setting, the reverse-engineer LLM performs  $M = \{5, 10, 20, 50\}$  rounds of interventions conditioned on the initial set of 10 observations ( $|\mathcal{O}| = 10$ ). In addition to GPT-4o, we report full results for Claude-3.5-Sonnet-20241022 (Anthropic, 2024), DeepSeek-R1 (Guo et al., 2025), Llama-3.3-70B-Instruct Grattafiori et al. (2024), GPT-5, and Claude 4 Sonnet in Appendix F, showing that even the strongest models still require active intervention to achieve high performance. We also show the reliability of using GPT-4o as a judge in Appendix I, comparing Cohen's scores with human annotations and with the latest LLM-as-judge results. We provide prompts for both intervention and hypothesis generation in Appendix E and other evaluation approaches in Appendix J.

#### 4.1 LLM Struggles to Utilize Observations Optimally

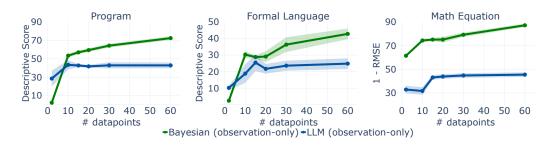


Figure 2: Observation-only results across three black-box types. We compare GPT-40 (blue) to Bayesian inference (green). The horizontal-axis represents the number of provided (x, y) pairs. We report 1 - RMSE for Math Equation and descriptive score for Program and Formal Language.

We first establish the reference performance achievable by the Bayesian model in each setting. These three settings were selected in part because they are all cases where previous work has defined inference algorithms that make it possible to approximate the posterior distribution over hypotheses as more observations becomes available (Rule et al., 2024; Yang & Piantadosi, 2022; Foster et al., 2019). As shown in Figure 2, the Bayesian models (green) consistently improve with the increased number of observations across all three tasks. On the other hand, while the LLM reverse-engineer (blue) starts off with higher performance for Program and Formal Language, potentially leveraging its prior knowledge, it peaks at 10 observations and struggles to use the extra observations thereby causing performance to plateau. We also calculate repeated measures ANOVAs (Girden, 1992) for each black-box type and found significant Model  $\times$  number of datapoints interactions for Program  $(F(5,10)=51.9,\ p<0.001)$ , Formal Language  $(F(5,10)=11.8,\ p=0.001)$ , and Math Equa-

tion (F(5,10)=8.7, p=0.002), showing that the Bayesian inference algorithms increasingly outperformed LLMs with additional observations. Details for the ANOVAs are in Appendix D.1.

#### 4.2 Intervention Is Crucial for the LLM to Refine Hypotheses

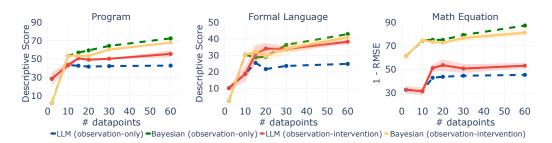


Figure 3: Observation-intervention results across three black-box types. Red: observations and interventions by GPT-4o. Yellow: taking the observation-intervention collected from GPT-4o as observations for the Bayesian inference algorithms. Dashed lines: observation-only reference for GPT-4o (blue) and Bayesian inference (green).

In Figure 3, we compare the performance of models with access to only the observations (dashed lines) against using the data that is actively collected through intervention (solid lines). We observe that enabling the LLM to actively intervene significantly improves performance (red) over observation-only (dashed blue). Through intervention, the LLM consistently improves as more data becomes available across all three black-box types, consistent with prior results on passive learning (Ostrovski et al., 2021). To assess the quality of the interventions, we provide the LLM-collected intervention data to the Bayesian model as observations, akin to the passive yoked data studied in Markant & Gureckis (2010; 2014). Our results indicate that while the interventions are beneficial to the LLM, they are not universally more informative, paralleling findings in human active learning (Markant & Gureckis, 2010; 2014). This gap was statistically significant, as shown by an ANOVA for each black box type: Program (F(5,10) = 23.9, p < 0.001), Formal Language (F(5,10) = 7.9, p = 0.003), and Math Equation (F(5,10) = 14.9, p < 0.001).

## 4.3 Identifying the Value of Generating the Intervention Data

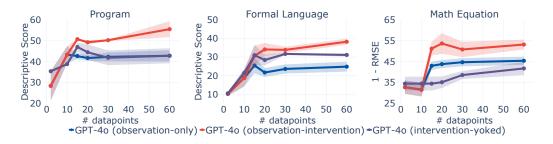


Figure 4: Comparing intervention-yoked results with observation-only and observation-intervention across three black-box types.

The improvement in performance produced by the interventions could have two sources: it could be that the resulting data are more informative, or that the process of generating interventions itself helps the model. To tease these apart, we adopt the passive-yoked design that Markant & Gureckis (2010; 2014) used to study human learning, where the data generated via active learning by one group of participants is presented to another group of participants as passive observations. In Figure 4, we compare GPT-40 across three conditions: **observation-only** (blue), **observation-intervention** (red), and **intervention-yoked** (purple) where GPT-40 only passively observes the interventional data without the verbalization and analysis that are used to construct such data. Results consistently show that the intervention-yoked setting leads to lower performance compared to

the observation-intervention setting across all three black-box types. This shows that active learning is more beneficial than passive-yoked learning in part because it allows the LLM to dynamically refine its hypothesis in response to its own interventions.

## 5 ANALYSIS

## 5.1 ESCAPING THE FAILURE MODES: OVERCOMPLICATION & OVERLOOKING

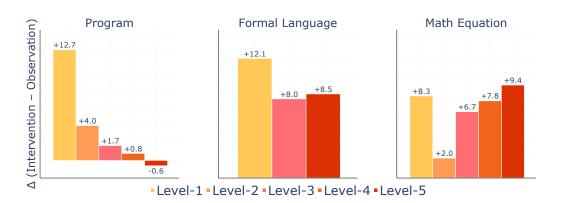


Figure 5: Descriptive scores for five different complexity levels. Averaged across three seeds for each of the three black-box types.

To understand how intervention improves LLM performance, we analyze common failures by sampling 20 failed examples (scoring below 2 out of 10 points) from the observation-only experiment, which were inspected by human experts. We provide more details in Appendix G.1. We identify two major failure modes: 1) *overcomplication*, where the LLM excessively interprets the data, resulting in unnecessarily complex hypotheses, and 2) *overlooking*, where the LLM inadequately leverages available information, leading to poorly reasoned hypotheses. We classified 20 randomly sampled examples for each black-box into the two failure modes or "Not Applicable" by human annotation. Results show that for Program the failures are predominantly from overcomplication (17 cases out of 20) whereas Math Equation contains more overlooking failures (16 cases out of 20). The failures are more evenly distributed for Formal Language, with 8 examples classified as overcomplication, 11 examples as overlooking, and 1 example as "Not Applicable". We provide examples for these failure mode in Appendix G.2.

Notably, we find that the impact of interventions on alleviating these two failure modes is contingent upon the complexity of the reverse-engineering task itself. For each of the three specific domains we study, we include a brief characterization of complexity in Appendix M. Within each domain, we observe that the complexity of the reverse-engineering problem instance characterized by  $f^*$  governs the extent to which interventions rectify overcomplication and overlooking. In Figure 5, we show that performance improvements from intervention on Program diminish as task complexity increases for black-box systems dominated by the *overcomplication* failure mode. In contrast, actively collected data is more beneficial when addressing challenging black-box instances dominated by the *overlooking* failure mode, such as Math Equation. For Formal Language, where both failure modes frequently occur, we observe consistent improvements across all complexity levels.

Case study. Figure 6 demonstrates how an LLM progressively updates its hypothesis through active interventions to ultimately reverse-engineer the underlying mechanism of a black-box system using a Formal Language black-box intervention example, where GPT-40 strategically designs subsequent queries to validate its current belief about the system. In contrast, under the observation-only scenario, the model remains trapped in identifying spurious patterns from passively observed data and lacks a meaningful way to assess its own uncertainty. Through active interventions, the LLM iteratively tests and revises its hypotheses after encountering failures, gradually reducing uncertainty and converging toward an accurate understanding of the black-box mechanism.

```
Black-Box: Formal Language Rule = A^n B^n C^(2n)

AAAABBBCCCCCCC is generated by the black box is generated by the black box is generated by the black box

(more observations)

Hypothesis: the black box generates strings that follow a pattern where [...]

- If not empty, it consists of a sequence of 'A's followed by 'B's followed by 'C's.

- The number of 'A's and 'B's can vary, but the number of 'C's tends to be larger.

Let's test strings with a different count of 'A's, 'B's, and 'C's to see if it is generated by the black box.

Test: AAABBBCCCC Black-box: cannot be generated

Test: AAABBBCCCC Black-box: cannot be generated

Test: AAABBBCCCC Black-box: cannot be generated

Test: AAABBBCCCC Black-box: can be generated

Test: AAABBCCCC Black-box: can be generated

Test: AAABBCCCC Black-box: can be generated

Test: AAABBCCCC Black-box: can be generated

Test: AAABBBCCCC Generated, indicating the number of A's and B's must be equal. [...]

- ABBCCC: Not generated, reinforcing the rule that C's must be exactly double the number of A's (or B's).

- AABBCCC: Not generated, reinforcing the rule that C's must be exactly double the number of A's (or B's).

- AABBCCC: Not generated, orifinming the pattern of equal A's and B's must be equal.

Hypothesis: the black box generates strings that follow these rules:

- If not empty, it consists of a sequence of 'A's followed by 'B's followed by 'C's.

- The number of 'A's and 'B's must be equal.

- The number of 'A's and 'B's must be equal.

- The number of 'A's and 'B's must be exactly double the number of A's (or 'B's).

Let's test a string that follows this pattern with a different count of A, B, and C.

Intervention Round N+1
```

Figure 6: Case study example. GPT-40 updates the hypothesis using intervention on Formal Language black box. Yellow: GPT-40 states the hypothesis based on the passive observations in round N and updates the hypothesis in round N+1. Red: test case. Teal: black box response.

#### 5.2 Intervention Strategies

Similar to how LLMs use chain-of-thought reasoning (Wei et al., 2022) to solve complex tasks, we allow the LLM to state its hypotheses and analyze the observations before constructing the query. We investigate how different reasoning strategies impact the effectiveness of interven-

Black Box	Intervention	Descriptive Intervention	Functional Intervention	Analyze-then- Query Intervention
Program	43.4	47.6	19.2	50.8
Formal Language	24.1	28.6	22.8	34.7
Math Equation	34.8	38.8	39.9	38.0

Table 1: Comparison of the four intervention strategies.

tion. We compare four strategies: 1) Intervention: no reasoning before constructing the query, 2) Descriptive Intervention: describing the current hypothesis about the black-box, 3) Functional Intervention: verbalize the black-box implementation as a Python program (Li et al., 2025; Luo et al., 2025), and 4) Analyze-then-Query: allowing the LLM to analyze data and state a hypothesis freely. Throughout our experiments, we allow the LLM to reason once every five queries.

As shown in Table 1, allowing the LLM to reason generally improves the effectiveness of intervention regardless of the strategy. However, the most effective intervention typically requires the LLM to carefully analyze past observations and strategically plan subsequent steps to acquire more informative data from the black-box. Interestingly, while structured reasoning in functional intervention (Li et al., 2025; Luo et al., 2025) is known to improve performance in formal reasoning tasks, it does not produce additional improvement in the context of reverse-engineering. This suggests that the LLM reverse-engineering abilities may differ from its formal reasoning capabilities.

## 5.3 Transferring to another LLM

We also examine whether interventional data actively collected by one LLM (GPT-4o) can effectively transfer and benefit another LLM (Llama-3.3-70B-Instruct). This is relevant to whether AI scientists can transfer their experiments and findings successfully to another AI scientist. Adopting a similar passive-yoked design, we compare three scenarios for Llama-3.3-70B-Instruct (Grattafiori et al., 2024): **observation-only**, **observation-intervention**, and **intervention-transfer**, where the interventional data is collected by GPT-4o. As shown in Figure 7, the intervention-transfer scenario achieves performance comparable to or slightly better than the observation-only baseline but consistently underperforms Llama's own intervention (observation-intervention). This suggests that while the intervention data from GPT-4o is informative, the effectiveness diminishes when transferred to a different LLM, showing that the benefit from intervention is model-specific.

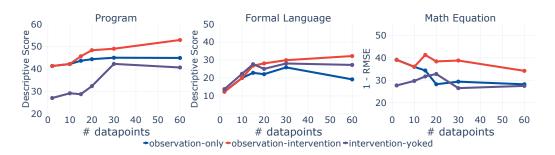


Figure 7: Intervention data transfer results. Red: Llama-3.3-70B-Instruct performing intervention. Blue: Llama-3.3-70B-Instruct using observations only. Purple: using interventional data from GPT-40 as observations for Llama-3.3-70B-Instruct.

# 6 LIMITATIONS AND FUTURE DIRECTIONS

In this paper, we have discussed the inabilities and failure modes of LLMs in reverse-engineering black-boxes. However, the three black-box types we studied represent only a narrow slice of possible tasks, even within controlled settings. A more comprehensive assessment will require expanding and scaling up the evaluation suite to probe LLMs' reverse-engineering abilities across a broader spectrum of scenarios. In addition, we have assumed idealized, noise-free black-boxes and fully trustworthy data—a condition that is rarely met in real scientific practices. An important next step is to relax this assumption and rigorously test LLM robustness in the presence of noise and uncertainty. As our paper discuss extensively on the failure modes of LLMs, we leave open the question: "How can we train LLMs to become effective reverse engineers?", which includes enhancing the LLM's ability to perform correct inference from passive observations and to conduct optimal experiments. In particular, what kinds of data and algorithms are needed to train such a model (for example, reinforcement learning using black-box environments), and can improvements in one domain generalize to broader scientific automation tasks? Finally, we have demonstrated that the actively acquired data by one LLM may not be useful for another LLM, pointing to the issue of experience transferability. Just as many major scientific advances have relied on effective human collaborations, so too may future automation of scientific discoveries depend on resolving this issue for LLM collaborations. Understanding and quantifying the impact of this limited transferability of knowledge may be crucial as multi-agent systems become prevalent, and it will be essential to design such systems with effective communication baked in.

## 7 Conclusion

In this paper, we identified and formalized the reverse-engineering problem as a core ability and prerequisite for performing a reliable scientific discovery. We showed that current LLMs still struggle to effectively leverage passive observations even on seemingly simple and controlled black-boxes. Allowing LLMs to actively collect intervention data improves performance, but still falls short of closing the gap with Bayesian inference, casting doubt on the promise of truly reliable AI scientists. Through extensive analysis, we identified issues such as overcomplication and overlooking and illustrate how intervention can mitigate such failures. Despite the effectiveness of intervention, our analysis revealed that the intervention data collected by LLMs were primarily beneficial to the models themselves, rather than being objectively informative or transferable to other models.

Altogether, our paper directly assesses the ability of LLMs to infer underlying causal structures and mechanisms through controlled reverse-engineering experiments. This capacity mirrors the fundamental scientific discovery process, which relies heavily on identifying hidden relationships and principles behind observed phenomena. Consequently, if an LLM cannot reliably reverse-engineer even simple or controlled systems, this raises concerns regarding its dependability in addressing more complex and ambiguous scientific challenges. Evaluating an LLM's reverse-engineering ability provides a concrete and principled way to assess its capacity for scientific reasoning, helping us understand whether such models possess the foundational skills required to function as dependable AI scientists.

## **8 ETHICS STATEMENT**

This work investigates how language models can reverse-engineer black-box systems in fully synthetic domains such as programs, formal languages, and mathematical equations. Our study does not involve human subjects, sensitive or proprietary data, or any real-world systems. While reverse-engineering methods in general could raise concerns if applied to sensitive settings, our research design deliberately avoids such cases by restricting all experiments to controlled, non-sensitive environments. We do not identify direct ethical issues beyond the standard considerations for computational research.

#### REFERENCES

- Jacob Abernethy, Kareem Amin, Michael Kearns, and Moez Draief. Large-Scale Bandit Problems and KWIK Learning. In *International Conference on Machine Learning*, pp. 588–596, 2013.
- Dhruv Agarwal, Bodhisattwa Prasad Majumder, Reece Adamson, Megha Chakravorty, Satvika Reddy Gavireddy, Aditya Parashar, Harshit Surana, Bhavana Dalvi Mishra, Andrew Mc-Callum, Ashish Sabharwal, et al. Open-ended scientific discovery via bayesian surprise. *arXiv* preprint arXiv:2507.00310, 2025.
- Pranjal Aggarwal, Seungone Kim, Jack Lanchantin, Sean Welleck, Jason Weston, Ilia Kulikov, and Swarnadeep Saha. Optimalthinkingbench: Evaluating over and underthinking in llms. *arXiv* preprint arXiv:2508.13141, 2025.
- Dana Angluin. Queries and Concept Learning. Machine Learning, 2:319–342, 1988.
- Dana Angluin and Carl H Smith. Inductive Inference: Theory and Methods. *ACM Computing Surveys (CSUR)*, 15(3):237–269, 1983.
- Anthropic. Claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet, 2024.
- Dilip Arumugam and Thomas L Griffiths. Toward Efficient Exploration by Large Language Model Agents. *arXiv preprint arXiv:2504.20997*, 2025.
- Richard Bellman. A Markovian Decision Process. *Journal of Mathematics and Mechanics*, pp. 679–684, 1957.
- Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Ronen I Brafman and Moshe Tennenholtz. R-MAX A General Polynomial Time Algorithm for Near-Optimal Reinforcement Learning. *Journal of Machine Learning Research*, 3:213–231, 2002.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. pp. 1877–1901, 2020.
- Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*, 2025.
- Kathryn Chaloner and Isabella Verdinelli. Bayesian Experimental Design: A Review. *Statistical Science*, pp. 273–304, 1995.

- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu,
   Mengfei Zhou, Zhuosheng Zhang, et al. Do Not Think that Much for 2+3=? On the Overthinking
   of o1-like LLMs. arXiv preprint arXiv:2412.21187, 2024.
  - Sahil Chopra, Michael Henry Tessler, and Noah D Goodman. The first crank of the cultural ratchet: Learning and transmitting concepts through language. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 41, 2019.
  - Julian Coda-Forno, Marcel Binz, Jane X Wang, and Eric Schulz. Cogbench: A Large Language Model Walks into A Psychology Lab. *arXiv preprint arXiv:2402.18225*, 2024.
  - Alejandro Cuadron, Dacheng Li, Wenjie Ma, Xingyao Wang, Yichuan Wang, Siyuan Zhuang, Shu Liu, Luis Gaspar Schroeder, Tian Xia, Huanzhi Mao, et al. The danger of overthinking: Examining the reasoning-action dilemma in agentic tasks. *arXiv* preprint arXiv:2502.08235, 2025.
  - MH DeGroot. Uncertainty, Information, and Sequential Experiments. *The Annals of Mathematical Statistics*, 33(2):404–419, 1962.
  - Marc Deisenroth and Carl E Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 465–472, 2011.
  - Armen Der Kiureghian and Ove Ditlevsen. Aleatory or Epistemic? Does it Matter? *Structural Safety*, 31(2):105–112, 2009.
  - Fabian Falck, Ziyu Wang, and Chris Holmes. Is in-context learning in large language models Bayesian? A martingale perspective. *arXiv* preprint arXiv:2406.00793, 2024.
  - Adam Foster, Martin Jankowiak, Elias Bingham, Paul Horsfall, Yee Whye Teh, Thomas Rainforth, and Noah Goodman. Variational Bayesian optimal experimental design. *Advances in Neural Information Processing Systems*, 32, 2019.
  - Michael C Frank. Baby steps in evaluating the capacities of large language models. *Nature Reviews Psychology*, 2(8):451–452, 2023.
  - Kanishk Gandhi, Michael Y Li, Lyle Goodyear, Louise Li, Aditi Bhaskar, Mohammed Zaman, and Noah D Goodman. BoxingGym: Benchmarking progress in automated experimental design and model discovery. *arXiv preprint arXiv:2501.01540*, 2025.
  - Yolanda Gil, Mark Greaves, James Hendler, and Haym Hirsh. Amplify scientific discovery with artificial intelligence. *Science*, 346(6206):171–172, 2014.
  - Ellen R Girden. ANOVA: Repeated measures. Number 84. Sage, 1992.
  - Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an AI co-scientist. arXiv preprint arXiv:2502.18864, 2025.
  - Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco

595

596

597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623 624

625

626 627

628

629 630

631

632

633 634

635

636

637

638

639

640 641

642

643

644

645

646

647

Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, and Zoe Papakipos. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783. Accessed: 2025-05-14.

- Thomas L Griffiths, Jian-Qiao Zhu, Erin Grant, and R Thomas McCoy. Bayes in the Age of Intelligent Machines. *Current Directions in Psychological Science*, 33(5):283–291, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Yancheng He, Shilong Li, Jiaheng Liu, Weixun Wang, Xingyuan Bu, Ge Zhang, Zhongyuan Peng, Zhaoxiang Zhang, Wenbo Su, and Bo Zheng. Can large language models detect errors in long chain-of-thought reasoning? *arXiv* preprint arXiv:2502.19361, 2025.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. GPT-40 system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Katie Kang, Amrith Setlur, Dibya Ghosh, Jacob Steinhardt, Claire Tomlin, Sergey Levine, and Aviral Kumar. What do learning dynamics reveal about generalization in llm reasoning? *arXiv* preprint arXiv:2411.07681, 2024.
- Alexander Ku, Declan Campbell, Xuechunzi Bai, Jiayi Geng, Ryan Liu, Raja Marjieh, R Thomas McCoy, Andrew Nam, Ilia Sucholutsky, Veniamin Veselovsky, et al. Using the tools of cognitive science to understand large language models at different levels of analysis. *arXiv* preprint *arXiv*:2503.13401, 2025.
- Salvatore La Torre, Parthasarathy Madhusudan, and Gennaro Parlato. A robust class of context-sensitive languages. In 22nd Annual IEEE Symposium on Logic in Computer Science (LICS 2007), pp. 161–170, 2007.

- Andrew Kyle Lampinen, Stephanie C Y Chan, Ishita Dasgupta, Andrew J Nam, and Jane X Wang. Passive learning of active causal strategies in agents and language models. *Advances in Neural Information Processing Systems*, 2023.
  - GG Landis JRKoch. The measurement of observer agreement for categorical data. *Biometrics*, 33 (1):159174, 1977.
  - Jia Li, Ge Li, Yongmin Li, and Zhi Jin. Structured chain-of-thought prompting for code generation. *ACM Transactions on Software Engineering and Methodology*, 34(2):1–23, 2025.
  - Lihong Li and Michael L Littman. Reducing Reinforcement Learning to KWIK Online Regression. *Annals of Mathematics and Artificial Intelligence*, 58:217–237, 2010.
  - Lihong Li, Michael L Littman, and Thomas J Walsh. Knows What It Knows: A Framework for Self-Aware Learning. In *Proceedings of the 25th International Conference on Machine Learning*, pp. 568–575, 2008.
  - Dennis V Lindley. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, 27(4):986–1005, 1956.
  - Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
  - Evan Z Liu, Aditi Raghunathan, Percy Liang, and Chelsea Finn. Decoupling Exploration and Exploitation for Meta-Reinforcement Learning Without Sacrifices. In *International Conference on Machine Learning*, pp. 6925–6935, 2021.
  - Ryan Liu, Jiayi Geng, Joshua C Peterson, Ilia Sucholutsky, and Thomas L Griffiths. Large language models assume people are more rational than we really are. *arXiv preprint arXiv:2406.17055*, 2024a.
  - Ryan Liu, Jiayi Geng, Addison J Wu, Ilia Sucholutsky, Tania Lombrozo, and Thomas L Griffiths. Mind your step (by step): Chain-of-thought can reduce performance on tasks where thinking makes humans worse. *arXiv preprint arXiv:2410.21333*, 2024b.
  - Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
  - Yijia Luo, Yulin Song, Xingyao Zhang, Jiaheng Liu, Weixun Wang, GengRu Chen, Wenbo Su, and Bo Zheng. Deconstructing long chain-of-thought: A structured reasoning optimization framework for long cot distillation. *arXiv* preprint arXiv:2503.16385, 2025.
  - Doug Markant and Todd Gureckis. Category learning through active sampling. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 32, 2010.
  - Douglas B Markant and Todd M Gureckis. Is it better to select or to receive? learning via active and passive hypothesis testing. *Journal of Experimental Psychology: General*, 143(1):94, 2014.
  - R Thomas McCoy and Thomas L Griffiths. Modeling rapid language learning by distilling bayesian priors into artificial neural networks. *arXiv* preprint arXiv:2305.14701, 2023.
  - R Thomas McCoy, Shunyu Yao, Dan Friedman, Mathew D Hardy, and Thomas L Griffiths. Embers of autoregression show how large language models are shaped by the problem they are trained to solve. *Proceedings of the National Academy of Sciences*, 121(41):e2322420121, 2024.
  - Ziqi Ni, Yahao Li, Kaijia Hu, Kunyuan Han, Ming Xu, Xingyu Chen, Fengqi Liu, Yicong Ye, and Shuxin Bai. Matpilot: an llm-enabled ai materials scientist under the framework of human-machine collaboration. *arXiv preprint arXiv:2411.08063*, 2024.
  - Charles O'Neill, Tirthankar Ghosal, Roberta Răileanu, Mike Walmsley, Thang Bui, Kevin Schawinski, and Ioana Ciucă. Sparks of science: Hypothesis generation using structured paper data. *arXiv* preprint arXiv:2504.12976, 2025.

- Ian Osband, Daniel Russo, and Benjamin Van Roy. (More) Efficient Reinforcement Learning via
   Posterior Sampling. Advances in Neural Information Processing Systems, 26:3003–3011, 2013.
  - Georg Ostrovski, Pablo Samuel Castro, and Will Dabney. The difficulty of passive learning in deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:23283–23295, 2021.
    - Martin L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 1994.
    - Tom Rainforth, Adam Foster, Desi R Ivanova, and Freddie Bickford Smith. Modern Bayesian Experimental Design. *Statistical Science*, 39(1):100–114, 2024.
    - Chandan K Reddy and Parshin Shojaee. Towards scientific discovery with generative ai: Progress, opportunities, and challenges. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 28601–28609, 2025.
    - Ronald L Rivest and Robert E Schapire. Diversity-Based Inference of Finite Automata. In 28th Annual Symposium on Foundations of Computer Science (SFCS 1987), pp. 78–87, 1987.
    - Ronald L Rivest and Robert E Schapire. Inference of Finite Automata Using Homing Sequences. In *Proceedings of the Twenty-First Annual ACM Symposium on Theory of Computing*, pp. 411–420, 1989.
    - Milena Rmus, Akshay K. Jagadish, Marvin Mathony, Tobias Ludwig, and Eric Schulz. Towards automation of cognitive modeling using large language models. *arXiv preprint arXiv:2502.00879*, 2025.
    - Joshua S Rule, Steven T Piantadosi, Andrew Cropper, Kevin Ellis, Maxwell Nye, and Joshua B Tenenbaum. Symbolic metaprogram search improves learning efficiency and explains rule learning in humans. *Nature Communications*, 15(1):6847, 2024.
    - Amin Sayedi, Morteza Zadimoghaddam, and Avrim Blum. Trading off Mistakes and Don't-Know Predictions. *Advances in Neural Information Processing Systems*, 23, 2010.
    - Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.
    - Burr Settles. Active Learning Literature Survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2009.
    - Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K Reddy. Llm-sr: Scientific equation discovery via programming with large language models. *International Conference on Learning Representations*, 2025.
    - Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs generate novel research ideas? A large-scale human study with 100+ NLP researchers. *arXiv preprint arXiv:2409.04109*, 2024.
    - Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To CoT or not to CoT? chain-of-thought helps mainly on math and symbolic reasoning. *arXiv preprint arXiv:2409.12183*, 2024.
    - Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, et al. Paperbench: Evaluating ai's ability to replicate ai research. *arXiv preprint arXiv:2504.01848*, 2025.
    - Alexander L Strehl and Michael L Littman. An Analysis of Model-based interval estimation for Markov Decision Processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
    - Malcolm JA Strens. A Bayesian framework for reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pp. 943–950, 2000.

- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Hanjie Chen, Xia Hu, et al. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.
- Richard S Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the Seventh International Conference on Machine Learning*, pp. 216–224, 1990.
- Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.
- Richard S Sutton and Andrew G Barto. Introduction to Reinforcement Learning. MIT Press, 1998.
- István Szita and Csaba Szepesvári. Agnostic KWIK learning and Efficient Approximate Reinforcement Learning. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 739–772, 2011.
- Sebastian B Thrun and Knut Möller. Active exploration in dynamic environments. *Advances in Neural Information Processing Systems*, 4, 1991.
- Thomas J Walsh, István Szita, Carlos Diuk, and Michael L Littman. Exploring Compact Reinforcement-Learning Representations with Linear Regression. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pp. 591–598, 2009.
- Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific Discovery in the Age of Artificial Intelligence. *Nature*, 620(7972):47–60, 2023.
- Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, et al. Thoughts are all over the place: On the underthinking of o1-like llms. *arXiv* preprint arXiv:2501.18585, 2025.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems, 35:24824–24837, 2022.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit Bayesian inference. In *International Conference on Learning Representations*, 2021.
- Yuan Yang and Steven T. Piantadosi. One model for the learning of language. *Proceedings of the National Academy of Sciences*, 119(5):e2021865119, 2022.
- Lance Ying, Katherine M Collins, Lionel Wong, Ilia Sucholutsky, Ryan Liu, Adrian Weller, Tianmin Shu, Thomas L Griffiths, and Joshua B Tenenbaum. On benchmarking human-like intelligence in machines. *arXiv preprint arXiv:2502.20502*, 2025.
- Cedegao E Zhang, Katherine M Collins, Lionel Wong, Mauricio Barba, Adrian Weller, and Joshua B Tenenbaum. People use fast, goal-directed simulation to reason about novel games. *arXiv* preprint arXiv:2407.14095, 2024.
- Jian-Qiao Zhu and Thomas L Griffiths. Eliciting the priors of large language models using iterated in-context learning. *arXiv preprint arXiv:2406.01860*, 2024a.
- Jian-Qiao Zhu and Thomas L Griffiths. Incoherent probability judgments in large language models. *arXiv preprint arXiv:2401.16646*, 2024b.

## A APPENDIX

#### B BLACK BOX DESIGNS

**Program** We used 100 list-mapping program instances from (Rule et al., 2024) to design the Program black-box API. Each black-box instance represents as a symbolic program defined in a domain-specific language (DSL). We implemented an interpreter pipeline that parses DSL expressions into abstract syntax trees and compiles them into executable Python code.

Each black-box supports two modes: observation (observation-only) and intervention (observation-intervention). In the observation mode, the black-box takes a random input list and returns the output produced by the underlying symbolic program, generating paired observational data:

input list 
$$\rightarrow$$
 program execution  $\rightarrow$  output list

In the intervention mode, the LLM queries an input or explicitly specifies an input-output pair. The black-box generates the output list or evaluates whether the given output matches the internally computed output and provides clear feedback:

$$Feedback = \begin{cases} \text{"output} \Rightarrow Correct", & \text{if the provided output matches the program output,} \\ \text{"output} \Rightarrow Incorrect", & \text{otherwise.} \end{cases}$$

**Formal Language** We followed (Yang & Piantadosi, 2022; McCoy & Griffiths, 2023) to implement a collection of 46 formal language instances to construct our formal language black-box, each instance being capable of generating strings according to specific symbolic rules  $(e.g.\ A^nB^n)$ . Each black-box instance behaves as an API from a generative model, operating in two modes: observation and intervention.

In the observation mode (observation-only), the black-box randomly produces valid strings from its underlying rule, explicitly labeling each as generated output, for example:

In the intervention mode (observation-intervention), the LLM submits a specific string query to the black-box, which evaluates whether the string complies with its rule. The black-box responds clearly, indicating either acceptance or rejection:

$$Response = \begin{cases} \text{``[string] is generated by the black-box''}, & \text{if the strings compile with the rule,} \\ \text{``[string] cannot be generated by the black-box''}, & \text{otherwise.} \end{cases}$$

To avoid generating infinite strings, we imposed a maximum character length of 64 for all single characters generated by the black-box.

**Math Equation** For the math equation, we implemented the CES utility model as the black-box, designing it as a generative model capable of generating observational data or responding to queries from an LLM. The utility function of CES is mathematically defined as:

$$U = \left(\sum_{i} a_{i} x_{i}^{r}\right)^{\frac{1}{r}},$$

where the weights  $a_i$  satisfy the constraint  $\sum_i a_i = 1$ , the parameter r controls the substitution elasticity, and  $x_i$  represents the quantities of goods in a basket.

CES black-box also provides two operational modes: observation (observation-only) and intervention (observation-intervention). In the observation mode, the black-box randomly samples two baskets (each a list of good quantities) and computes their utilities using the CES formulation. It then returns the preference outcome indicating which basket is preferred based on

higher utility:

$$\mbox{Preference} = \begin{cases} \mbox{Basket1}, & U(\mbox{Basket1}) > U(\mbox{Basket2}), \\ \mbox{Basket2}, & U(\mbox{Basket1}) < U(\mbox{Basket2}), \\ \mbox{equal utility}, & U(\mbox{Basket1}) = U(\mbox{Basket2}). \end{cases}$$

In the intervention mode, an external model explicitly queries the black-box by specifying two baskets. In addition, the external model can also provide an estimated preference. The CES black-box internally evaluates the utilities based on the specified baskets and returns the actual preference outcome or feedback indicating whether the provided estimate was "correct" or "incorrect".

# C BAYESIAN MODELS AS THE 'OPTIMAL' REFERENCE

We employ Bayesian models as an oracle for optimal reverse-engineering against which we may assess the capabilities of LLMs. Unlike LLMs, Bayesian models explicitly perform probabilistic inference within a clearly defined hypothesis space, systematically updating posterior beliefs using the Bayes rule to identify the underlying mechanism that best explain observed data. Under the critical assumption that the true underlying rule resides within this hypothesis space (that is, the standard assumption of a well-specified prior), Bayesian models serve as an optimal reference standard in our experimental setting. We hypothesize that LLMs, when provided only with passive observational data, are unable to effectively utilize available information due to their inherent reliance on prior knowledge, resulting in significantly lower performance compared to the Bayesian optimal standard. However, allowing LLMs to actively intervene and collect data can substantially reduce the performance gap. For each of the three black-box systems evaluated, we replicated the Bayesian models from their original studies, adapting them to closely match our experimental conditions. Specifically, we provide Bayesian models with observed data generated by our black-box systems as an ideal reference. We also provide Bayesian models with the actively collected data from LLMs intervention to assess the informativeness of the data gathered by LLMs. To ensure rigorous comparability, we applied identical evaluation methodologies to both the Bayesian models and LLMs.

**Program** We used the Bayesian inference approach from Rule et al. (2024) to establish an optimal reference for list-mapping program black-box. Specifically, we utilized their MetaProgram Learner, which performs Bayesian inference over symbolic metaprograms that generate target programs from observed data.

Given observational data D, consisting of input-output pairs generated by symbolic programs, the MPL computes the posterior distribution over candidate hypotheses (metaprograms) H according to the Bayes rule:

$$P(H \mid D) \propto P(D \mid H) \cdot P(H)$$
.

The prior distribution P(H) integrates two complementary sources of simplicity bias: the meta-program prior  $P_{\mathcal{M}}(H)$  and the induced program prior  $P_{\mathcal{P}}(\widetilde{H})$ . This combined prior is defined as:

$$P(H) \propto \exp\left(\frac{\ln P_{\mathcal{M}}(H) + \ln P_{\mathcal{P}}(\widetilde{H})}{2}\right),$$

where  $\widetilde{H}$  denotes the program compiled from the metaprogram H.

The likelihood  $P(D \mid H)$  measures the consistency of a meta-program H with the observational data provided, incorporating a noise model to accommodate minor discrepancies between the model predictions and observations.

**Formal Language** We adopted the Bayesian inference approach from (Yang & Piantadosi, 2022) as an optimal reference model to determine the theoretical upper bound on the learnability of formal language rules from the observations generated by our black-boxes or from the interventions queried by LLM. Specifically, we provided strings generated by our formal language black-boxes

as observational data to the Bayesian model, which then inferred the underlying symbolic grammar rules.

Just as before, the Bayesian inference framework defines the posterior distribution over candidate hypotheses conditioned on observed data using Bayes' rule:

$$P(H \mid D) \propto P(D \mid H) P(H),$$

where H represents a candidate hypothesis (grammar or probabilistic program), D represents the observed string data generated by the black-box, P(H) represents the prior probability reflecting initial beliefs about the simplicity and plausibility of hypotheses, and  $P(D \mid H)$  denotes the likelihood of observing data D given hypothesis H.

The Bayesian model uses a structured prior P(H), assigning higher probabilities to simpler, more concise grammars or symbolic programs. As observational data increases, Bayesian updating systematically refines prior beliefs into posterior distributions, enhancing the probability assigned to grammars that best explain the data. Formally, each new observed string updates the posterior, shifting probability mass toward hypotheses consistent with the cumulative dataset. By leveraging this Bayesian inference mechanism, we quantify the upper bound of the learnability of the observations, thus providing a rigorous baseline to evaluate LLM's effectiveness in utilizing the same observational data.

**Math Equation** To infer the parameters of the CES utility model from the observations provided, we followed (Foster et al., 2019) by employing a Bayesian inference approach explicitly conditioned on these observations. Bayesian inference integrates observed data with prior beliefs, updating these beliefs into posterior distributions to progressively improve parameter estimates. Initially, we specified prior distributions for the model parameters:

```
\rho \sim \text{Beta}(\rho_0, \rho_1),
\alpha \sim \text{Dirichlet}(\alpha_{\text{conc}}),
\text{slope} \sim \text{LogNormal}(\text{slope}_{\mu}, \text{slope}_{\sigma}).
```

Given pairs of consumption bundles  $(d_1, d_2)$  and the corresponding observed user preferences y, the Bayesian framework models these preferences probabilistically through a censored sigmoid-normal likelihood:

```
y \sim \text{CensoredSigmoidNormal} \left( \text{slope} \cdot (U(d_1) - U(d_2)), \text{ slope} \cdot \text{obs\_sd} \cdot (1 + ||d_1 - d_2||_2) \right),
```

where  $U(d_1) - U(d_2)$  denotes the utility difference between the two bundles. Here, "censored" refers to applying a sigmoid function to latent utility values and then truncating the results to the observed preference interval (e.g., [0,1]), ensuring that responses remain within these limits.

The posterior distributions are updated via Bayes' theorem by explicitly integrating observational data:

```
p(\rho, \alpha, \text{slope} \mid y, d) \propto p(y \mid \rho, \alpha, \text{slope}, d) p(\rho, \alpha, \text{slope}),
```

where  $p(\rho, \alpha, \text{slope})$  represents prior distributions and  $p(y \mid \rho, \alpha, \text{slope}, d)$  represents the likelihood function given the observations.

While some sources prefer uppercase probability notation such as  $P(H \mid D)$ , this paper adopts lowercase notation (p) consistently for both probability densities and random variables throughout.

Parameter estimation was performed via variational inference (Blei et al., 2017), iteratively optimizing the evidence lower bound (ELBO), defined as:

```
\mathsf{ELBO}(\phi) = \mathbb{E}_{q_{\phi}} \left[ \log p(y \mid \rho, \alpha, \mathsf{slope}, d) \right] - D_{\mathsf{KL}} \left( q_{\phi}(\rho, \alpha, \mathsf{slope}) \parallel p(\rho, \alpha, \mathsf{slope}) \right),
```

where  $q_{\phi}$  denotes the variational posterior distribution used to approximate the true posterior distribution.

Thus, as additional observational data are incorporated, Bayesian inference continually updates prior beliefs into posterior distributions, systematically refining parameter estimates toward their true underlying values.

## D STATISTICAL SIGNIFICANT TESTS

#### D.1 REPEATED-MEASURES ANOVA

 To statistically evaluate the interaction between models (Bayesian vs. LLM) and steps, we calculated the repeated-measures ANOVAs. Each black-box instance involved multiple repeated measurements corresponding to different steps. Letting  $Y_{ijk}$  represent the performance score for subject i, models j (Bayesian or LLM), and step k, the repeated-measures ANOVA model can be expressed as:

$$Y_{ijk} = \mu + S_i + M_j + T_k + (M \times T)_{jk} + \epsilon_{ijk}$$

where  $\mu$  is the mean in all measurements,  $S_i$  represents the random effect of the subjects (individual seeds),  $M_j$  denotes the main effect of the model,  $T_k$  is the main effect of steps,  $(M \times T)_{jk}$  is the interaction between the model and the step, and  $\epsilon_{ijk}$  represents residual error.

The ANOVA decomposes the total variance into these distinct sources. Specifically, the significance of the interaction of the Step Method  $\times$  was determined by calculating the corresponding F-statistic:

$$F = \frac{MS_{(M \times T)}}{MS_{error}}$$

where  $MS_{(M\times T)}$  is the mean square for the Method  $\times$  Step interaction, and  $MS_{error}$  is the residual mean square. Significance was assessed using an F-distribution with numerator degrees of freedom equal to (J-1)(K-1), where J is the number of method levels and K is the number of steps, and denominator degrees of freedom equal to (I-1)(K-1), where I is the number of subjects.

#### E PROMPTS

## E.1 Intervention prompt

In this task, you are given a `black box'' and need to determine its inner workings. {black box information}
You will have a series of turns to interact with the black box. On each turn, you can either gather more information or test your hypothesis. To gather more information, {query instruction}, and obtain a result.

To test your hypothesis, {test instruction}. All the information gathered across all the turns is used to reverse engineer the black box. Throughout the process, you can decide whether the gathered information is sufficient to correctly identify the workings of the black box, in which case you can stop. Otherwise, you need to continue the interaction. Concretely, you can perform one of the following actions at each turn: 1) query, 2) test, or 3) stop.

Provide a \*thorough reasoning\* before performing the action. Leverage the past observations to design your next query and make your hypothesis as accurate as possible. Below is the format for each action.

```
Query:
```query
List[int]
```

Test:

```
1026
      ```test
1027
      List[int]
1028
      List[int]
1029
1030
      Stop:
1031
      ```stop
1032
1033
1034
      Note that you should only perform one of the actions above with
1035
      one input example in your response.
1036
1037
      Below are your past observations of the black box:
1038
     {observations}
     Response:
1039
1040
1041
      E.2 EVALUATION PROMPTS
1042
      Program (judge):
1043
1044
      In this task you will be given a ground truth program and
1045
      pseudocode that you need to evaluate. You will output a score for
1046
      the quality of the pseudocode based on a set of assessment
1047
      criteria.
1048
1049
      Below is the ground truth program:
1050
      {ground_truth}
1051
      Evaluate the quality of the following pseudocode:
1052
     {response}
1053
1054
      Score the above pseudocode against the ground truth program based
1055
      on the following criteria (total 10 points):
1056
      1. Does the provided pseudocode correctly specify the
1057
      implementation of the ground truth program and manipulate the
1058
      variables in the same way? Ignore the programming language
1059
      difference. [5 point]
1060
      2. Does the provided pseudocode specify the implementation in the
      most simple and straightforward way without extra unused parts
1061
      (Occam's Razor principle) [5 point]
1062
1063
      Explain your judgement and return the final score with the type
1064
      float and following the format below:
1065
      ```judgement
1066
      YOUR JUDGEMENT HERE
1067
      ···score
1068
      YOUR SCORE HERE
1069
1070
1071
      Response:
1072
1073
      Formal Language (judge):
1074
1075
      In this task, you will be given a ground truth formal language and
1076
      a proposed rule describing that formal language, which you need to
1077
      evaluate for quality. You will then output a score based on a set
```

Below is the ground truth formal language:

of assessment criteria.

```
1080
      {ground_truth}
1081
      Evaluate the quality of the following formal language rule:
1083
      Score the above formal language rule against the ground truth
      formal language based on the following criteria (total: 10
1085
      points):
      1. Does the provided rule correctly generate the examples given in
1087
      the ground truth? Your score is determined by how many examples
1088
      are correctly generated out of the total number of examples. [3
1089
      points]
1090
      2. Does the provided rule correctly reverse engineer the ground
1091
      truth formal language? [5 point]
1092
      3. Is the provided rule in the most simple and straightforward way
1093
      without extra unused parts (Occam's Razor principle)? Note: If the
      provided rule is completely incorrect, you should give 0 point for
1094
     this criterion. [2 point]
1095
1096
      Explain your judgement and return the final score with the type
1097
      float and following the format below:
1098
      ```judgement
1099
      YOUR JUDGEMENT HERE
1100
1101
      ```score
1102
      YOUR SCORE HERE
1103
1104
      Response:
1105
1106
      Math Equation (judge):
1107
1108
      In this task, you are provided with a ground truth CES utility
1109
      function and a CES utility function predicted by a model.
1110
      Your task is to evaluate the quality of the predicted utility
1111
      function based on a set of assessment criteria and output a score.
1112
1113
      The ground truth utility takes this form:
1114
      U(\mathbb{x}) = \left(\frac{i=1}\right^n a_i \cdot \left(\frac{i}{x}\right)
1115
      x_i^{{\\text{{rho}}}}\\right)^{{1/\\text{{rho}}}}
1116
1117
      The utility depends on the following parameters:
1118
      1. a_i: float rounded to the first decimal point and should sum up
1119
      to 1. (Note that there will be multiple a i's.)
      2. rho: float rounded to the first decimal point.
1121
      Below is the information about the ground truth utility function:
1122
     {ground_truth}
1123
1124
      Evaluate the quality of the following predicted the parameters of
1125
     the utility function:
1126
      {response}
1127
1128
      Score the predicted utility function against the ground truth
      using the following criteria (total 10 points):
1129
      1. Is the predicted utility function has a correct rho? [2 points]
1130
      2. Compare the predicted utility function to the ground truth, how
1131
     many a_i's are correct (order matters)? This will give us an
1132
      accuracy percentage. The score for this bullet should be the
1133
      accuracy percentage times the total allocated 6 points [6 points]
```

```
1134
      3. In the predicted utility function, do the unknown parameters
1135
      a_i sum up to 1 and do the number of a_i's match the number of
1136
      goods? [1 point]
1137
      4. Does the predicted utility function express the function in a
1138
      simple and straightforward way without any unnecessary elements
      (adhering to the Occam's Razor principle)? [1 point]
1139
1140
      Explain your judgement and return the final score with the type
1141
      float and following the format below:
1142
      ```judgement
1143
      YOUR JUDGEMENT HERE
1144
1145
1146
      ```score
1147
      YOUR SCORE HERE
1148
1149
      Response:
1150
1151
      Descriptive Evaluation:
1152
1153
      In this task, you are given a "black box" and need to determine
1154
      its inner workings.
1155
      {black box information}
1156
1157
      Below are some past observations from the black box:
      {observations}
1158
1159
      Your task is to reverse engineer the rule underlying {more
1160
      detailed instructions} in the following format:
1161
       ``Rule
1162
      YOUR RULE HERE
1163
1164
1165
      Response:
1166
      Function Implicit Evaluation:
1167
1168
      In this task, you are given a "black box" and need to determine
1169
      its inner workings.
1170
      {black box information}
1171
      Below are some past observations from the black box:
1172
      {observations}
1173
1174
      {More detailed instructions}
1175
      Output your generated string in the following format:
1176
      ```output
1177
      YOUR RESPONSE HERE
1178
1179
1180
      Response:
1181
1182
         REVERSE ENGINEERING ABILITIES ACROSS DIFFERENT CATEGORIES OF
1183
         LLMs
1184
1185
```

In Figure 8, we report the results of observation-only and observation + intervention across different LLMs: Llama-3.3-70B-Instruct, Claude-3.5-Sonnet, DeepSeek-R1, GPT-5 and Claude 4 Sonnet. Across nearly all black-box types and models, actively intervening with iteratively refining hypothe-

1186

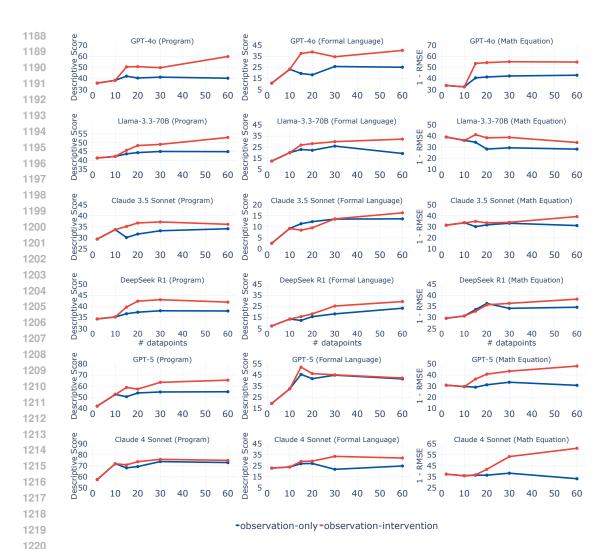


Figure 8: Results of reverse engineering abilities across different categories of LLMs. We report Llama-3.3-70B-Instruct, Claude 3.5 Sonnet, Deepseek R1, GPT-5, and Claude 4 Sonnet on Program, Formal Language and Math Equation.

ses consistently enhances models' understanding of the underlying black-box dynamics. In particular, we show that DeepSeek R1, Claude 4 Sonnet and GPT-5, utilizing Long CoT reasoning, have the potential to continuously extract informative knowledge even from passive learning scenarios. This detailed and long reasoning allows the model to explore various potential hypotheses. However, despite these advantages, frontier reasoning LLMs do not significantly outperform models without explicit reasoning (e.g., GPT-40, Llama-3.3-70B-Instruct, Claude 3.5 Sonnet) in reverse-engineering tasks, except for Claude 4 on program. This finding highlights the inherent limitations of current reasoning steps for existing LLMs.

## G COMMON FAILURE MODES

#### G.1 HUMAN ANNOTATION

To systematically analyze LLM's failure modes, we defined an LLM reverse-engineering attempt as a failure if its descriptive score was below 2 out of 10, according to our descriptive evaluation rubric. For each black-box type, we randomly selected 20 representative failure cases from the observation-only setting. We have two human experts independently reviewed these examples, categorizing each

case based on the nature of the error. Any disagreements were resolved through discussion. Finally, human annotators identified two common failure modes: overcomplication and overlooking.

#### G.2 Overcomplication & Overlooking Examples

Across the three black-box types, we find that overcomplication is a common failure mode, particularly in the Program, while overlooking most often occurs in Math Equation. For Formal Language, both overcomplication and overlooking are observed when LLMs fail at reverse engineering. In Tables 2,3 4 and 5, we show the failure examples for Program (overcomplication), Formal Language (overcomplication & overlooking) and Math Equation (overlooking).

## H COMPLEXITY CATEGORIZATION

We rank the complexity level from 1-5. Each black-box type includes multiple instances of varying task complexity.

**Program.** The complexity level is determined based on the number of operations, which ranges from 1-12. Instances with fewer than 2 operations are classified as complexity level 1 (complexity-1), those with fewer than 4 operations as complexity-2, fewer than 6 operations as complexity-3, and fewer than 8 operations as complexity-4. Due to the limited number of remaining examples, all others are grouped into the highest complexity level (complexity-5).

**Formal Language.** Instead of using five complexity levels, we divided the Formal Language instances into three levels, drawing on insights from (La Torre et al., 2007). Specifically, we categorized regular language instances as complexity-1 black-boxes, context-free languages as complexity-3, and context-sensitive languages as complexity-5.

**Math Equation.** We categorize complexity levels according to the number of goods involved, ranging from 2 to 6. Specifically, instances with 2 goods are labeled as complexity - 1, 3 goods as complexity - 2, and so on, with instances involving 6 goods classified as the highest complexity level, complexity - 5.

## I RELIABILITY AND ACCURACY OF USING GPT-40 AS A JUDGE

The use of LLM-as-Judge has been a common practice to evaluate model generation and GPT-4 level models have been shown to match or exceed human annotation in quality (??) for evaluating generated text. In our experiment settings, the LLM judge takes a set of rubrics that sum to a total of 10 points, and the description of the black-box instance to score the model response description of the black-box instance to score the model response. Our implementation further removes the potential vagueness by adding rubrics to evaluate the correctness in a fine-grained manner. The description of the black-box instances are also non-ambiguous to the model as we provide the context in which they need to be interpreted. We show GPT-40's reliability as a judge by computing Cohen's kappa between GPT-40 and (i) thinking LLMs (OpenAI o3 and Claude-4-Sonnet) and (ii) human annotations. We randomly sample 30 examples (10 for each black-box type) and collect annotations to calculate the Weighted Cohen Kappa score (for ordinal rating). We obtain an overall Weighted Cohen Kappa score of 0.773 for Human, 0.752 for Claude 4, and 0.734 for o3. All the results indicate substantial agreement (Landis JRKoch, 1977) and show the reliability and accuracy of using GPT-40 as a judge.

# J EVALUATION OF THE REVERSE-ENGINEERING ABILITY

Unlike typical tasks used to benchmark LLMs, such as solving math problems or question answering which are commonly evaluated using accuracies, the reverse-engineering ability is less straightforward. One can assess how well the black-box  $f^*$  is recovered by an LLM by: 1) descriptive evaluation where the LLM verbalizes the hypothesis to compare to the ground truth and 2) functional evaluation which captures how well the LLM emulates the black-box's input-output dynamics and

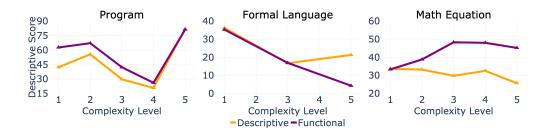


Figure 9: Comparison of descriptive evaluation (yellow) and functional evaluation (purple) across black-box complexity levels.

generalizes to unseen examples (Kang et al., 2024). In functional evaluation, the LLM directly predicts the response conditioned on the test query and the past observations and compute accuracy  $\text{Acc} = \frac{1}{M} \sum_{i=1}^{M} \mathbb{I}[y_i^{\text{test}} = \mathcal{M}(x_i^{\text{test}}, \mathcal{O})]$ , without generating the black-box implementation, akin to in-context learning (Brown et al., 2020). As shown in Figure 9, descriptive and functional evaluation trends align for Program across complexity levels. However, we also observe discrepancies of trends between the two evaluations for Formal Language (complexity level 3 to 5) and Math Equation (complexity level 1 to 3), demonstrating that the evaluation protocol and the *format* of the model output may capture different strengths and weaknesses of the model. For Program, we used the original samples from the black box of the list mapping program as test cases (Rule et al., 2024) and ensured that none of these input—output pairs were included in the observations. For Formal Language and Math Equation, we use our deterministic black-box randomly sample 20 test cases per black-box instance.

## K ANOTHER BLACK-BOX TYPE: BOARD GAME

#### K.1 Black-Box Design

We design a connect-n board game (2  $\times$  2 to 9  $\times$  9) variant following (Zhang et al., 2024). The black-box is defined by the rules that dictate the winning condition of the game (e.g., "Win by connecting 3 stones in a column."). The LLM can query with a board state and an action, and the black-box responds with the new board state and a game status (win/lose/draw/ongoing). In our black-box design, every game instance exposes two modes—observation (observation-only) and intervention (observation-intervention) —and uses the symbols X and 0 to mark the moves of the two players.

**Game definition.** For a given instance, let the board be a  $r \times c$  grid and let  $\langle r_{\text{win}}, c_{\text{win}}, d_{\text{win}} \rangle$  denote the required number of consecutive marks needed to win horizontally, vertically, and diagonally, respectively. During play the black-box tracks the current board state B, the active player  $p \in \{X, O\}$ , and whether the game has ended.

In observation mode, an external LLM supplies an *initial* board (or leaves it empty). The blackbox generates the following as the outputs:

- the round number,
- · the updated board,
- whose move it was last,
- the current game status (ongoing, draw, PlayerX\_wins, etc.).

If the move ends the game, the record also names the winner.

In intervention, the LLM needs to specify (i) additional pieces to place on the board, (ii) the candidate action it wishes the black-box to take, and (iii) optionally, a predicted follow-up board. The black-box returns the same structured record as in observation mode. If the LLM also proposed a prediction of the next state, the black-box confirms it ("Correct") or explains why it is invalid. For Board Game, we do not have a Bayesian model as the optimal reference for the comparison.

#### K.2 GPT-40 RESULTS

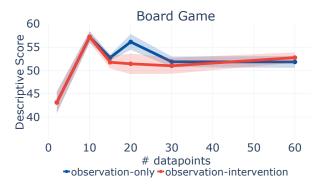


Figure 10: Observation-only and observation-intervention results for Board Game.

In Figure 10, we do not observe the same trends seen in Programs, Formal Language, and Math Equation black-box types. For Board Game, actively collected data does not improve the reverse-engineering performance of the model, indicating that the data gathered is not even significantly informative for the LLM itself. We hypothesize that this is because, to query the black-box, the LLM must (1) generate a board state, (2) propose a next move, and (3) predict the resulting board state, requiring a multi-step reasoning process. These compounded requirements make it challenging for the LLM to probe edge cases or effectively reduce uncertainty about the black-box. This result further highlights a key limitation of current LLMs: When the information signal from the black-box is sparse, actively collected data remain of limited utility.

## L FUNCTIONAL EVALUATION DETAILS

For Program, we used the original samples from the black box of the list mapping program as test cases (Rule et al., 2024) and ensured that none of these input—output pairs were included in the observations. For Formal Language and Math Equation, we use our deterministic black-box randomly sample 20 test cases per black-box instance.

## M COMPLEXITY CATEGORIZATION

We rank the complexity level from 1-5. Each black-box type includes multiple instances of varying task complexity.

**Program.** The complexity level is determined based on the number of operations, which ranges from 1-12. Instances with fewer than 2 operations are classified as complexity level 1 (complexity-1), those with fewer than 4 operations as complexity-2, fewer than 6 operations as complexity-3, and fewer than 8 operations as complexity-4. Due to the limited number of remaining examples, all others are grouped into the highest complexity level (complexity-5).

**Formal Language.** Instead of using five complexity levels, we divided the Formal Language instances into three levels, drawing on insights from (La Torre et al., 2007). Specifically, we categorized regular language instances as complexity-1 black-boxes, context-free languages as complexity-3, and context-sensitive languages as complexity-5.

**Math Equation.** We categorize complexity levels according to the number of goods involved, ranging from 2 to 6. Specifically, instances with 2 goods are labeled as complexity - 1, 3 goods as complexity - 2, and so on, with instances involving 6 goods classified as the highest complexity level, complexity - 5.

1454 1455

```
1405
1406
        Black-box instance: (lambda (singleton (third $0)))
1407
1408
        Observations:
                           Input: [97, 53, 5, 33, 65, 62, 51]; Output: [5]
1409
                           Input: [61, 45, 74, 27, 64]; Output: [74]
1410
                           Input: [36, 17, 96]; Output: [96]
1411
                           Input: [79, 32]; Output: invalid input
1412
                           Input: [90, 77, 18, 39, 12, 93, 9, 87, 42]; Output: [18]
1413
                           Input: [71, 12, 45, 55, 40, 78, 81, 26]; Output: [45]
                           Input: [61, 56, 66, 33, 7, 70, 1, 11, 92]; Output: [66]
1414
1415
                           Input: [90, 100, 85, 80, 0, 78, 63]; Output: [85]
1416
                           Input: [31, 93, 41, 90, 8, 24]; Output: [41]
1417
                           Input: [28, 30, 18, 69, 57, 11, 10, 40, 65, 62]; Output: [18]
                           Input: [38, 70]; Output: invalid input
1418
                           Input: [90, 15, 70, 42, 69]; Output: [70]
1419
                           Input: [77, 70, 75, 36]; Output: [75]
1420
1421
                           Input: [11, 76, 49, 40, 73, 30, 37, 23]; Output: [49]
                           Input: [23, 4, 78, 84]; Output: [78]
1422
                           {More observations}
1424
1425
                           Input: [60, 8, 11, 86, 96]; Output: [11]
1426
                           Input: [19, 4, 10]; Output: [10]
1427
                           Input: [87, 50, 90, 67, 35, 66, 30, 27, 86]; Output: [90]
1428
                           Input: [53, 74, 35, 57, 63, 84, 82, 89, 45, 10]; Output: [35]
1429
                           Input: [78, 14, 62, 75, 80, 42]; Output: [62]
1430
                           Input: [31, 2, 93, 34]; Output: [93]
1431
                           Input: [90, 28]; Output: invalid input
1432
                           Input: [21, 42, 54, 7, 12, 100]; Output: [54]
1433
                           Input: [89, 28, 5]; Output: [5]
1434
                           Input: [81, 68, 77, 87, 9, 3, 15, 81, 24, 77]; Output: [77]
1435
                           Input: [15, 50, 11, 47, 14, 4, 77, 2, 24, 23]; Output: [11]
1436
                           Input: [61, 26]; Output: invalid input
1437
                           Input: [86]; Output: invalid input
1438
                           Input: [79, 12, 33, 8, 28, 9, 82]; Output: [33]
1439
                           Input: [44, 55, 23, 7, 64]; Output: [23]
1440
1441
1442
                               FUNCTION black_box(input_list):
1443
                               IF length of input_list; 3:
                               RETURN "invalid input"
1444
                               ELSE:
1445
                               SORT input_list in ascending order
1446
                               RETURN the third element in the sorted input_list
1447
        Model response:
1448
        (Overcomplication)
1449
                               Explanation:
                               The black box function takes a list of integers as input. If the list contains fewer than
1450
                               3 elements, it returns "invalid input". Otherwise, it sorts the list in ascending order
1451
                               and returns the third smallest element (which is the third element in the sorted
1452
                               list).
1453
```

Table 2: Overcomplication example of Program. Red: LLM comes up with a more complex rule of black-box instance.

Black-box instance: AnBmAnBm  $(A^n B^m A^n B^m)$ **Observations:** AAAAAAAAAAA is generated by the black box AABBBBBBAABBBBBB is generated by the black box AAAAAAAAAAAAAAABB is generated by the black box AAAAAAAABBBBBAAAAAAAABBBBB is generated by the black box AAAAABBBBBBBBAAAAABBBBBBBB is generated by the black box AAAAAAABBBBBBAAAAAAAABBBBBB is generated by the black box {More observations} AAABBBBAAABBBB is generated by the black box AAAABBAAAABB is generated by the black box AAAAAAAAAAAAAAA is generated by the black box AAAAAABBBBBBAAAAAABBBBBB is generated by the black box AAAAAABBBBBAAAAAABBBBB is generated by the black box AAAAAAAAAAAAAA is generated by the black box AAAAAAAAABBBBBBBBBAAAAAAAAAABBBBBBBB is generated by the black box AAAAAAAAAABBAAAAAAABB is generated by the black box AAAAABBBBBBBBAAAAABBBBBBBB is generated by the black box AAAABBBBAAAABBBB is generated by the black box Model response: The formal language generates strings composed of two segments: a sequence of 'A's fol-lowed by a sequence of 'B's. The number of 'A's in the first segment must be equal to (Overcomplication) the number of 'B's in the second segment, or the number of 'A's can be zero, resulting in a string of only 'B's. Additionally, the language allows for strings composed entirely of 'A's or entirely of 'B's. The sequences of 'A's and 'B's can be repeated multiple times, as long as the rule of equal counts between 'A's and 'B's in each segment is maintained. 

Table 3: Overcomplication example for Formal Language. Red: The LLM predicts a rule that is more complex than necessary for this black-box instance. In this example, it incorrectly reverse the number of As to be equal to the number of Bs.

1512 1513 1514 1515 1516 1517 Black-box instance: XXX (three copies of the same string) 1518 **Observations:** BBABBABBABBABBABBABBABA is generated by the black box 1519 BABBABBABBABBABBAB is generated by the black box 1520 ABBABBABB is generated by the black box 1521 ABAAAABABAAAABABAAAAB is generated by the black box 1522 AAA is generated by the black box AAAAAA is generated by the black box BBB is generated by the black box 1525 ABAABAABAABAABA is generated by the black box ABABAB is generated by the black box ABAABAABA is generated by the black box {More observations} 1531 BBBBAABABABBBBAABABABBBBAABABA is generated by the black box 1532 ABABAB is generated by the black box 1533 AAABAABAABAABAABAAB is generated by the black box 1534 ABABAB is generated by the black box 1535 1536 ABBBABBBBBBBB is generated by the black box 1537 BABABA is generated by the black box 1538 BBABAABBABBABAABBABBABAABBA is generated by the black box 1539 AAAAAA is generated by the black box 1540 1541 Model response: The black box generates strings that are composed of repeated patterns or segments. Each 1542 (Overlooking) string can be broken down into one or more segments, where each segment is a repetition of 1543 a smaller pattern. The patterns can be composed of the letters 'A' and 'B', and the repetition of these patterns can vary in length. Additionally, the black box can generate an empty The key characteristics of the generated strings are: 1547 1. The strings can be empty. 1548 2. The strings can be composed of repeated patterns of 'A' and 'B'. 1549 3. The patterns themselves can vary in complexity and length. 4. The repetition of these patterns can be consistent or vary within the string. 1550 1551 Examples of patterns include: 1552 - Single letters repeated (e.g., "AAA", "BBB"). 1553 - Simple alternating patterns (e.g., "ABABAB", "BABABA"). - More complex patterns with varying repetition (e.g., "ABBABBABB", "BABBABBAB-1554 BABBABBAB"). 1555 1556 The rule can be estimated as: The black box generates strings by repeating pat-1557 terns of 'A' and 'B', where the patterns can be simple or complex, and the repetition can be consistent or vary within the string. 1559

Table 4: Overlooking example of Formal Language. LLM estimates a rule of the black-box instance without utilizing the observed data to come up with the hypothesis.

1560

1561

#### 1566 Black-box instance: alpha=[0.1,0.1,0.4,0.4]; rho=0.4 1567 1568 Observations: Basket1: [87.79948556831374, 39.184851241810016, 1569 22.551465649464433, 42.84633678341351]; 1570 Basket2: [66.20150946367784, 80.43281524767865, 68.02177295663886, 8.112415033884117]; 1571 Preference: Basket1 1572 Basket1: [20.086807088533046, 79.79230548488334, 1573 8.817190035426963, 90.22168650956948]; 1574 Basket2: [90.40796323190197, 70.51080246059688, 1575 69.96510352125082, 48.476064999077785]; 1576 Preference: Basket2 Basket1: [8.003392409139845, 35.23615132704936, 1578 76.3120116750855, 59.18156984237496]; 1579 Basket2: [31.278783818535036, 18.047644469275518, 1580 21.855445800494433, 53.88293391636068]; 1581 Preference: Basket1 1582 {More observations} 1584 Basket1: [40.80654893485903, 14.795240036848078, 1586 37.35074210377093, 79.44298024697972]; 1587 [51.30239830060435, 50.0604427932334, 21.59530896237239, 36.88204913376821]; 1588 Preference: Basket1 Basket1: [75.18968775093136, 42.30088453163497, 1590 18.388385810882728, 95.92185012513752]; 1591 Basket2: [93.37477559149177, 86.19812315165629, 1592 49.438229718409715, 98.68538582332832]; Preference: Basket2 1594 Basket1: [70.30309661179918, 94.48475506394722, 1595 78.28267517818938, 70.90519145824096]; Basket2: [63.648862012782175, 99.80230448540196, 1597 61.33979160499423, 78.38413565179029]; 1598 Preference: Basket1 Model response: To reverse engineer the CES utility function, we need to determine the parameters $a_i$ and $\rho$ (Overlooking) that best fit the observed preferences between baskets. The CES utility function is given by: 1601 1602 $U(\mathbf{x}) = \left(\sum_{i=1}^{n} a_i \cdot x_i^{\rho}\right)^{1/\rho}$ 1603 1605 where $a_i$ are the weights for each good and $\rho$ is the elasticity parameter. Steps to Determine Parameters: 1606 Analyze Preferences: For each pair of baskets, the basket with the higher utility is preferred. 1609 • Estimate $\rho$ : The parameter $\rho$ affects the degree of substitutability between goods. 1610 $\rho=1$ is perfect substitutes, $\rho=0$ is Cobb-Douglas, $\rho<0$ is complementarity. 1611 • Estimate $a_i$ : The weights $a_i$ sum to 1 and reflect the importance of each good. 1612 • Iterative Fitting: Use optimization to adjust $a_i$ and $\rho$ to minimize preference viola-1613 1614 Estimation: After analyzing the data and testing various combinations, the following 1615 provides a good fit for the observed preferences: Best Fit Utility Function: 1616 1617 $U(\mathbf{x}) = (0.3x_1^{0.5} + 0.2x_2^{0.5} + 0.3x_3^{0.5} + 0.2x_4^{0.5})^2$ 1618 1619 This is derived by fitting the parameters to minimize preference violations. The weights $a_i$ reflect the relative importance of each good, and $\rho = 0.5$ is selected to match substitutability.