

ACTIVELY LEARNING HORN ENVELOPES FROM LLMs

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose a new strategy for extracting Horn rules from Large Language Models (LLMs). Our approach is based on Angluin’s classical exact learning framework where a learner actively learns a target formula in Horn logic by posing *membership* and *equivalence* queries. While membership queries naturally fit into the question-answering setup of LLMs, equivalence queries are more challenging. Previous works have simulated equivalence queries in ways that often lead to a large number of uninformative queries, making such simulations prohibitively expensive for modern LLMs. Here, we propose a new adaptive prompting strategy for posing equivalence queries, making use of the generative capabilities of modern LLMs to directly elicit counterexamples. We consider a case study where we learn rules describing gender-occupation relationships from the LLM. We evaluate our approach on the number of queries asked, the alignment of the extracted rules with historical data, and the consistency across runs. We show that our approach is able to extract Horn rules aligned with historical data with high confidence while requiring orders of magnitude fewer queries than previous methods.

1 INTRODUCTION

Large language models (LLMs) have in recent years displayed remarkable capabilities when interacting with users and answering queries. However, as LLMs are trained on large amounts of data, the information that is retained after training, the knowledge embedded in the model, becomes obscured. An important area of research in recent years has therefore been focused on assessing or extracting the knowledge acquired by LLMs (Roberts et al., 2020; Petroni et al., 2019; He et al., 2025). One can see the task of extracting knowledge from an LLM as an *active learning* task, where the learner interacts with a teacher by posing queries. The basic idea is that, by asking the right queries based on the answers received, the learner can navigate the space of possibilities and obtain the desired knowledge from the teacher more efficiently.

The exact learning framework by Angluin (1987) formalizes this mode of learning in a mathematical way. The most studied communication protocol between the learner and the teacher, also called an *oracle*, consists of *membership* and *equivalence queries*. In set theory notation, a membership query asks ‘Is $x \in \mathcal{T}$?’ where \mathcal{T} , known as the *target*, is a set that represents what the learner wants to learn. In an *equivalence query*, the learner sends its idea of what the target is, the *hypothesis* \mathcal{H} , and asks if \mathcal{H} is equivalent to \mathcal{T} . In set theory notation, ‘Is $\mathcal{H} = \mathcal{T}$?’. The answer is ‘yes’ if they are equivalent or ‘no’ together with a *counterexample* illustrating the difference, that is, a counterexample $x \in \mathcal{H} \oplus \mathcal{T}$, where \oplus is the symmetric difference. We say that the counterexample is *negative* if it models the hypothesis but not the target, and *positive* if it models the target but not the hypothesis.

In this work, we implement the task of extracting rules from LLMs as a learning task within Angluin’s exact learning framework. The membership query asks whether a statement is possible in the “view” of the LLM, to which the LLM should answer ‘yes’ or ‘no’ (e.g., ‘is it possible for astronauts to be born in Antarctica?’), and we formalize these expressions using propositional Horn logic (e.g. $\neg(\text{astronaut} \wedge \text{Antarctica})$). We collect the answers of the LLM to form a *hypothesis*, a Horn formula that expresses the knowledge extracted so far. Membership queries naturally fit into the question-answering setup of LLMs but answering equivalence queries is more complex in this setting. The equivalence query requires the teacher to understand the given hypothesis, compare with the target representing its “view” of the world, and if it is not equivalent, return a counterexample illustrating where the hypothesis and the target differ.

054 [Blum et al. \(2024\)](#) extracted Horn expressions from BERT-based models ([Devlin et al., 2018](#); [Liu](#)
055 [et al., 2019](#)) using random sampling to simulate equivalence queries. The idea behind simulating
056 equivalence queries this way is to ask the model to classify randomly generated examples as positive
057 (‘yes’) or negative (‘no’) with membership queries. If the learner finds an example x that is classified
058 differently by the model than expected from our hypothesis, then the learner can proceed as if the
059 teacher had replied ‘no’ to an equivalence query and returned x as counterexample. However, this
060 simulation strategy has several drawbacks, as random generation may create several uninformative
061 queries before finding a counterexample. We avoid simulating equivalence queries with random
062 sampling by *directly* querying for a counterexample to a given hypothesis. Considering the generative
063 capability of LLMs, we explore their ability to parse the hypothesis and provide a valid counterex-
064 ample if the hypothesis differs from the target. Although parsing logical information from LLM
065 responses has been previously explored ([Ye et al., 2023](#)), our work is, to the best of our knowledge,
066 the first to explore this with the purpose of generating counterexamples to equivalence queries. Our
067 algorithmic strategy is based on the algorithm by [Blum et al. \(2024\)](#), which is a non-trivial adaptation
068 of the algorithm for exact learning Horn expressions by [Angluin et al. \(1992\)](#). The adapted algorithm
069 returns a *Horn envelope*, a Horn expression that is the closest approximation of the true target. To
070 evaluate our approach, we consider a case study on gender-occupation relationships in LLMs.

071 The main contributions of this work are as follows: (i) an implementation of the algorithm for
072 learning Horn envelopes [Blum et al. \(2024\)](#) adapted to learn from generative LLMs; (ii) an adaptive
073 prompt generation strategy, based on the current working hypothesis built using the responses of the
074 LLM, that implements, not simulates, equivalence queries by directly asking the LLM to generate
075 counterexamples; (iii) an experimental analysis of gender-occupation relationship rules extracted
076 from LLMs, where we assess the impact of model type and model size of the rules extracted and
077 analyze the consistency and confidence of extracted rules.

078 In the following section, we describe related works on using machine learning (ML) models as
079 teachers and as reasoners. In Section [3](#) we provide basic notions and notation used in this paper. In
080 Section [4](#), we discuss some challenges in extracting rules (in the format of Horn expressions) from
081 LLMs, in particular, we describe our prompting strategy for dealing with equivalence queries. We
082 describe our experiments in Section [5](#) and the evaluation criteria and results in Section [6](#). Finally, we
083 conclude in Section [7](#).

084 2 RELATED WORKS

085 **ML models as teachers** Previous works have explored the idea of using a machine learning model as
086 a teacher within Angluin’s exact learning framework to extract information from the model. [Weiss](#)
087 [et al. \(2024\)](#) employed the exact learning framework to extract automata from recurrent neural
088 networks (RNNs). They used an algorithm designed to pose membership and equivalence queries.
089 Since RNNs cannot naturally answer equivalence queries, these queries were simulated in two ways:
090 by random sampling and by using a heuristic to find counterexamples. As previously mentioned,
091 [Blum et al. \(2024\)](#) extracted Horn expressions from BERT-based models and it is this work that
092 comes closest to ours. They conducted a case study on gender bias in occupations and extracted Horn
093 expressions manifesting biases in the BERT-based models. The same study has also been taken as
094 basis in the work by [Ozaki et al. \(2025\)](#), which extracts decision trees, also indicating the presence of
095 occupational-based gender biases in these models.

096 **ML models as reasoners** There have been many survey papers on reasoning in LLMs ([Mondorf &](#)
097 [Plank, 2024](#); [Huang & Chang, 2023](#); [Sun et al., 2024](#)) which explore reasoning behaviours in LLMs
098 and highlights the interest in the field. Previous works have employed many tools in aiding LLMs
099 in reasoning tasks. As mentioned [Ye et al. \(2023\)](#) used LLMs to parse logical information from
100 the query but they leveraged an outside theorem prover to derive the final answer. Our work uses
101 an algorithm to derive our hypothesis but still relies on the model to parse our hypothesis and use
102 reasoning to provide a counterexample. [Brown et al. \(2020\)](#) explored how few-shot prompting could
103 aid reasoning in larger models and this is the method used in our work.

3 PRELIMINARIES

We introduce the notation and basic relevant notions about propositional logic, Horn envelopes, and the exact learning framework.

Propositional Logic, Horn Envelopes Let \mathcal{P} be a set of propositional symbols. A *literal* is an element p of \mathcal{P} or its negation, denoted $\neg p$. A literal is *positive* if it is in \mathcal{P} , otherwise it is *negative*. A *clause* is a disjunction of literals, in symbols, $l_1 \vee \dots \vee l_n$, with $n \in \mathbb{N}$ and each l_i a (positive or negative) literal. A *Horn clause* is a clause with at most one positive literal. A *propositional expression* (in normal form) is a set of clauses. A *Horn expression* is a propositional expression that has only Horn clauses. The semantics of propositional expressions is given by *interpretations*. An interpretation is a subset of \mathcal{P} . An interpretation \mathcal{I} *satisfies*: a positive literal $p \in \mathcal{P}$ iff $p \in \mathcal{I}$; a negative literal $\neg p$ iff $p \notin \mathcal{I}$; a clause iff \mathcal{I} satisfies at least one of its literals; and a propositional expression iff it satisfies all of its clauses. Given a propositional expression ϕ , we write $\mathcal{I} \models \phi$ if \mathcal{I} satisfies ϕ . We define $\text{interpretations}(\phi)$ as $\{\mathcal{I} \mid \mathcal{I} \subseteq \mathcal{P}, \mathcal{I} \models \phi\}$. Given propositional expressions ϕ, ψ , we say that ϕ *entails* ψ , written $\phi \models \psi$, iff $\text{interpretations}(\phi) \subseteq \text{interpretations}(\psi)$. Also, ϕ, ψ are *logically equivalent*, written $\phi \equiv \psi$, iff $\text{interpretations}(\phi) = \text{interpretations}(\psi)$. The *closure under intersection* of a set of interpretations M is the set of all interpretations that can be obtained as the intersection of interpretations in M . Let $\text{Pow}(S)$ denote the power set of a set S .

Proposition 3.1 (Dechter & Pearl, 1992) *A set of interpretations $M \subseteq \text{Pow}(\mathcal{P})$ is closed under intersection iff $M = \text{interpretations}(\phi)$ for some Horn expression ϕ .*

By Proposition 3.1 for every set of interpretations closed under intersections, there is a unique (up to logical equivalence) Horn expression that represents it. Given an arbitrary propositional expression ϕ , the *Horn envelope* of ϕ is the Horn expression that represents the closure under intersection of the models of ϕ . The Horn envelope of ϕ , denoted $\text{env}(\phi)$, can be seen as the closest approximation of ϕ to a Horn expression. Syntactically, this Horn expression is assumed to be the Duquenne-Guigues basis (Guigues & Duquenne, 1986), which has the minimum number of rules. We often write clauses in the format of a *rule* $P \rightarrow Q$, where P is the conjunction of negative literals (or \top if none is negative) and Q is the disjunction of positive literals (or \perp if none is positive). E.g., $\neg p_1 \vee \neg p_2 \vee \dots \vee \neg p_n$ becomes $(p_1 \wedge p_2 \wedge \dots \wedge p_n) \rightarrow \perp$ and $\neg p_1 \vee \neg p_2 \vee \dots \vee \neg p_n \vee q$ turns into $(p_1 \wedge p_2 \wedge \dots \wedge p_n) \rightarrow q$. Intuitively, one can read $P \rightarrow Q$ as ‘if all propositional symbols in P hold then at least one of those in Q needs to hold’. In case there are no positive literals, that is, $P \rightarrow \perp$, then we read it as ‘a world where everything in P holds is not possible’.

Exact Learning We briefly present the exact learning framework (Angluin, 1987), adapted to the case of learning the Horn envelope of a propositional expression from interpretations. Membership and equivalence queries in this case are defined as follows. Let \mathcal{T} be a propositional expression. We call \mathcal{T} the *target*, meaning that \mathcal{T} represents the knowledge of the teacher. A *membership query for \mathcal{T}* takes as input an interpretation \mathcal{I} and returns ‘yes’ if \mathcal{I} satisfies \mathcal{T} and ‘no’ otherwise. In symbols, $\text{MQ}_{\mathcal{T}}(\mathcal{I}) = \text{yes}$ if $\mathcal{I} \in \text{interpretations}(\mathcal{T})$, otherwise $\text{MQ}_{\mathcal{T}}(\mathcal{I}) = \text{no}$. A *Horn equivalence query for \mathcal{T}* takes as input a propositional expression \mathcal{H} —the hypothesis—and returns ‘yes’ if the Horn envelopes of \mathcal{T} and \mathcal{H} are equivalent, otherwise, it returns ‘no’ and a counterexample in the symmetric difference of the sets of interpretations satisfying \mathcal{T} and \mathcal{H} . In symbols, $\text{EQ}_{\mathcal{T}}^{\text{Horn}}(\mathcal{H}) = \text{yes}$ if $\text{env}(\mathcal{T}) \equiv \text{env}(\mathcal{H})$, otherwise, ‘no’ and an element $\mathcal{I} \in \text{interpretations}(\mathcal{T}) \oplus \text{interpretations}(\mathcal{H})$. If the target and the hypothesis are Horn expressions then the notion of a Horn equivalence query coincides with the notion of equivalence query by Angluin et al. (1992), however, we cannot guarantee or expect the target to be a Horn expression when treating a machine learning model as the teacher. An algorithm *exactly learns* Horn expressions (with membership and Horn equivalence queries) if it always terminates returning a hypothesis equivalent to the Horn envelope of the target expression.

4 ACTIVELY LEARNING FROM LLMs

Extracting rules from LLMs is a challenging task as LLMs can struggle with logical reasoning (Mondorf & Plank, 2024) and are sensitive to changes in the prompt (Chen et al., 2024). For us to actively learn from the LLMs we need the model to be able to parse our hypothesis, compare with the target hypothesis, and reason about counterexamples in the symmetric difference between the hypothesis and the target. We describe the strategies employed to address the following challenges.

1. Translation: the format of the membership and equivalence queries needs to be adapted from interpretations and logical formulas into expressions in natural language, requiring translation from logic to natural language.
2. Format: we may ask the LLM to reply in a certain format, but there is no guarantee that the LLM will reply in the requested format, requiring us to validate responses to ensure a correct translation between logic expressions and natural language.
3. Validation: even if the LLM replies in the correct format, the reply may be incorrect, inconsistent with previous responses, or even change if the same query is posed multiple times.

Each of these challenges are carefully considered in our work. We start by explaining how we address the first challenge in the paragraph on *membership queries*. We describe how we address the two challenges within the paragraph on *equivalence queries* and in the paragraph on *Validating EQ responses*. Only once all challenges have been addressed are we able to translate the response from the model back into a logical expression for the algorithm.

Learning Horn Envelopes [Angluin et al. \(1992\)](#) presented a classical algorithm that exactly learns any target Horn expression in polynomial time. However, the polynomial bound is under the assumption that the target is indeed a Horn expression. This can no longer be assumed when we take an LLM as a teacher, and the classical algorithm is not guaranteed to terminate without the assumption. The main theoretical contribution of the work by [Blum et al. \(2024\)](#) was to propose an algorithm, called Horn Envelope, that handles the case where the unknown target is not (equivalent to) a Horn expression. The algorithm can be seen in the Appendix as Algorithm 2.

Membership Queries The role of the membership query is to allow us to uncover the model’s implicit knowledge and iteratively refine the rules that form the hypothesis of the learner. E.g., asking whether an interpretation that contains the symbols `Antarctica`, `before_1900` satisfies the “view” of the LLM, represented by \mathcal{T} , would correspond to the query $\text{MQ}_{\mathcal{T}}(\{\text{Antarctica}, \text{before_1900}\})$, which would translate to the prompt: “Has ‘A person born before 1900 in Antarctica’ been possible in the real world? Reply with ‘It is possible’ if yes and ‘it is NOT possible’ if no. Do not specify which values are different”¹.

Equivalence Queries The role of the equivalence query is to check if our hypothesis is equivalent to the target, and, if not, receive counterexamples that iteratively brings us closer to the target. If our current hypothesis was that no one had been born in Antarctica before 1900, then the query of if this hypothesis was equivalent to the the LLM’s “view” of the world, represented by \mathcal{T} , would correspond with $\text{EQ}_{\mathcal{T}}^{\text{Horn}}(\{\neg(\text{Antarctica} \wedge \text{before_1900})\})$, which would translate to appending: “The updated hypothesis: ‘A person born in any continent, in any time period, who is any occupation, could be any gender.’ describes the real world accurately, except A person born in Antarctica before the year 1900 is not possible. Please provide another counterexample to my hypothesis if possible” to the base equivalence query prompt. The complete few-shot equivalence prompt can be seen in appendix A.2.

Example Let the current hypothesis be empty, meaning that anything is possible. If the model contains the information that there were no astronauts born before 1900 then it could provide us with this information as a counterexample represented by the interpretation $\mathcal{I} = \{\text{astronaut}, \text{before_1900}\}$. Since this interpretation does not satisfy the target, we want to construct a hypothesis that is also not satisfied by this interpretation. The Horn Envelope algorithm does this by adding the Horn clause $\neg(\text{astronaut} \wedge \text{before_1900})$ to its hypothesis. In the next iteration, if the model responds to $\text{EQ}_{\mathcal{T}}^{\text{Horn}}(\{\neg(\text{astronaut} \wedge \text{before_1900})\})$ with the interpretation $\mathcal{J} = \{\text{astronaut}, \text{Antarctica}\}$ as a counterexample, since counterexamples are from the symmetric difference of the interpretations satisfying the target and the hypothesis, we know that this counterexample does not satisfy the target (there were no astronauts born in Antarctica) because it satisfies our hypothesis. The Horn Envelope algorithm tests all intersections of received negative counterexamples using membership queries. In this case, the intersection of \mathcal{I} and \mathcal{J} is $\mathcal{I} \cap \mathcal{J} = \{\text{astronaut}\}$. The algorithm would then ask a

¹The inclusion of “Do not specify which values are different” in the prompt was motivated by the models tendency to explain its answer in unpredictable ways. These explanations made parsing the response more difficult.

membership query with this intersection $\text{MQ}_{\mathcal{T}}(\{\text{astronaut}\})$. Suppose that the answer is ‘Yes’. Then the algorithm considers the interpretation $\{\text{astronaut}\}$ as one that satisfies the target. Then, the rule $\neg(\text{astronaut})$ is *not* added to the hypothesis. So, the counterexample $\mathcal{I} = \{\text{astronaut}, \text{Antarctica}\}$ is unchanged and the algorithm adds the rule $\neg(\text{astronaut} \wedge \text{Antarctica})$ to its hypothesis. The updated hypothesis would be: $\{\neg(\text{astronaut} \wedge \text{Antarctica}) \wedge \neg(\text{astronaut} \wedge \text{before}_{.1900})\}$

Testing intersections helps to create more informative counterexamples. When the model is no longer able to give a counterexample, the hypothesis is considered to be equivalent to the Horn envelope of the target, which is the closest Horn approximation of the propositional expression representing the relationship between the variables in the model. We note that rules of the form $(\text{astronaut} \wedge \text{before}_{.1930}) \rightarrow \text{male}$ (expressing that ‘astronauts born before 1930 are male’) can also be generated by the algorithm. This can happen when an interpretation that is not satisfied by the target is contained in an interpretation that is satisfied.

Validating EQ responses While LLMs are capable of responding to queries they do not always respond with a valid counterexample following the requested format. To ensure the validity of the responses to the equivalence queries, we check and handle a number of ways the responses might be invalid. The full list of steps can be seen in Algorithm 1, and a taxonomy of the validation steps can be seen in Table 1. If the validation checks discovers an error, then the EQ prompt is modified to reflect the error and the LLMs is queried again for a counterexample. In the case where the LLM is unable to understand the error description, to avoid repetitive loops, a threshold is set such that if the LLM makes too many validation errors in a row², then we consider the model unable to provide a valid counterexample and terminate the run.

Algorithm 1: Equivalence Query Response Validation

input : Equivalence Oracle $\text{EQ}_{\mathcal{T}}^{\text{Horn}}(\mathcal{H})$, hypothesis \mathcal{H} , list of previous counterexamples L
output : Validated *counterexample* or hypothesis \mathcal{H}

- 1 prompt \leftarrow translate \mathcal{H} to natural language and ask for counterexamples
- 2 **while** *the error threshold is not reached* **do**
 - 3 response \leftarrow query $\text{EQ}_{\mathcal{T}}^{\text{Horn}}(\mathcal{H})$ with prompt
 - 4 **if** *done thinking* **then**
 - 5 **if** *response is to terminate* **then**
 - 6 **return** $(\mathcal{H}, \text{termination flag})$
 - 7 **if** *response can be parsed* **then**
 - 8 potential counterexample \leftarrow parse response
 - 9 **if** *potential counterexample is not a duplicate* **then**
 - 10 **if** *potential counterexample is logically valid* **then**
 - 11 **return** validated counterexample
 - 12 **append** response with description of error to the prompt
- 13 **return** $(\mathcal{H}, \text{non-termination flag})$

5 EXPERIMENTS: GENDER-OCCUPATION RELATIONSHIPS IN LLMs

All experiments were performed on a High-Performance Computing (HPC) cluster with 1 GPU-accelerated node per experiment. The specific GPUs were either a NVIDIA A100 or a NVIDIA RTX3090, limited to 90GB of memory per run.

To allow a meaningful comparison with earlier research, we consider the same occupational gender bias scenario presented in Blum et al. (2024). The authors generate sentences using a template with a

²In our experiments the threshold of consecutive errors was set to 10.

Table 1: Taxonomy of response validation

| Validation step | Definition |
|--------------------------|---|
| <i>done thinking</i> (4) | The Deepseek reasoning model used in the experiments prepares its responses with a set of <think> tags. If there is no closing tag then the model did not complete its thinking, not providing a valid response. |
| <i>terminate</i> (5) | If the termination phrase is in the response then we return the current hypothesis \mathcal{H} . |
| <i>parse</i> (7) | If the response follows the requested response format, only contains values from the allowed set, and the response does not contain multiple counterexamples, then we parse the potential counterexample from natural language to logical language. |
| <i>duplicate</i> (9) | If the potential counterexample is not a duplicate of a previously given counterexample. |
| <i>valid</i> (10) | If the counterexample is a valid counterexample, that is, if the counterexample is in the symmetric difference between the hypothesis and target. |

limited set of values for each variable³ seen in Table 2. We recreate their experiments⁴ using their method, referred to in this paper as the *sampling method*, and compare with our own *direct EQ* approach.

Table 2: Set of allowed values

| Variable | Values |
|--------------------|---|
| Continent | ‘Africa’, ‘Americas’, ‘Asia’, ‘Europe’, ‘North America’, ‘Oceania’, ‘South America’, ‘Australia’, or ‘Eurasia’. |
| Time Period | ‘before 1875’, ‘between 1875 and 1925’, ‘between 1925 and 1951’, ‘between 1951 and 1970’, or ‘after 1970’ |
| Occupation | ‘fashion designer’, ‘nurse’, ‘dancer’, ‘priest’, ‘footballer’, ‘banker’, ‘singer’, ‘lawyer’, ‘mathematician’, or ‘diplomat’ |
| Gender | ‘woman’, or ‘man’ |

Direct EQ In the direct EQ approach we generate a prompt that communicates our current hypothesis and elicits counterexamples from the model directly. To avoid influencing the model by giving examples of possible and/or impossible combinations, we employed a zero-shot strategy for the first counterexample communicating the zero-hypothesis (everything is possible), instructions describing the format of the replies, and restrictions on the domain found in Table 2. We ask the model to provide us with a counterexample in the following form: “A person born in <continent>, <time period>, who is a <occupation> CAN/CANNOT be a <gender>.” Once a counterexample has been given, both the updated hypothesis and all previous counterexamples are communicated to the model in a few-shot setting.

Models For the experiments with the sampling method we used the RoBERTa-base and RoBERTa-large models (Liu et al., 2019) also used by Blum et al. (2024) as well as the instruction tuned Mistral-7B (Jiang et al., 2023) generative model. For the experiments with the direct EQ approach we used the instruction tuned Mistral-7B and Mistral-24B models as well as the DeepseekR1-8B reasoning model distilled from Llama (DeepSeek-AI, 2025).

6 EVALUATION AND RESULTS

We evaluate the resulting hypotheses by applying the following three criteria.

Number of Queries Prompting generative LLMs can be resource intensive. The sampling approach used by Blum et al. (2024) may generate many examples before a valid counterexample is found. A

³They also include ‘unknown value’, relevant to the Horn algorithm but unnecessary for our work as counterexamples generated with our approach simply omits variables that are not set. E.g. “A person born time period in continent.” if the gender is not set.

⁴We made the change to run their experiments to termination instead of limiting the experiment to 200 equivalence queries as was done by Blum et al. (2024).

324 formula for the upper bound of the number of queries needed is given by the Probably Approximately
 325 Correct (PAC) framework (Valiant, 1984; Angluin, 1987). The formula is $\frac{\ln(\frac{|H|}{\delta})}{\epsilon}$ where $|H|$ is the
 326 size of the hypothesis space (that is, the set of all possible hypotheses), ϵ represents the error, and δ
 327 represents the confidence.⁵ This upper bound can be quite large, so reducing the number of queries
 328 needed to generate the hypothesis is an important aspect of the evaluation.
 329

330 **Consistency of Hypotheses** Consistency of the hypotheses across multiple runs helps build con-
 331 fidence that the extracted hypothesis is equivalent to the target. To determine whether the model-
 332 generated hypotheses really represent the target, we evaluate the variance in the similarity of the
 333 extracted hypotheses. We evaluate the similarity by directly comparing clauses of the hypotheses
 334 for equality. Logically equivalent rules are treated as equal (e.g. $\neg(p \wedge q)$ and $\neg(q \wedge p)$ are logically
 335 equivalent and treated as equal in our calculation). Larger intersections between pairs of hypotheses
 336 with smaller variance indicate more consistent models.

337 As our baseline, we consider the expected size of the intersection of extracted rules between each
 338 pair of runs. Let S_1, S_2 denote their sets of rules, which are subsets of the set R of all possible
 339 rules. Assume that S_1, S_2 are random subsets of R . The expected size of the intersection $S_1 \cap S_2$
 340 given the sizes of S_1 and S_2 is

$$341 \quad E(|S_1 \cap S_2| \mid |S_1|, |S_2|) = \frac{|S_1| \cdot |S_2|}{|R|} \quad (1)$$

342 We compare the expected size of the intersection with the actual size of the intersection to determine
 343 whether the intersection could be attributed to randomness.
 344

345 **Correlation of rules with historical data** Wikidata⁶ contains publicly available data about his-
 346 torical persons including their occupations. We compare how the set of rules extracted from the
 347 LLMs corresponds with the data extracted from Wikidata. We should be aware that Wikidata contains
 348 historical data, which in some cases reflect historical gender biases with respect to occupation. Some
 349 examples of such biases can be seen in Figure 5 in the appendix.
 350

351 To study the correspondence between the rules extracted and real world information present in large
 352 datasets such as Wikidata, we calculate the *confidence* of the rule commonly used in association rule
 353 learning (Agrawal et al., 1993). The confidence is calculated with the following equation.
 354

$$355 \quad \text{conf}(X \rightarrow Y) = \frac{\text{supp}(X \cap Y)}{\text{supp}(X)} \quad (2)$$

356 The confidence of a rule is calculated by counting the number of people in Wikidata that support
 357 both the antecedent and the consequent of the rule, divided by the number of people that support the
 358 antecedent.⁷ Since the gender in our experiments is a binary variable limited to male or female, we
 359 consider that rules of the form $\neg(\text{continent} \wedge \text{time_period} \wedge \text{occupation} \wedge \text{woman})$ equivalent to
 360 $(\text{continent} \wedge \text{time_period} \wedge \text{occupation}) \rightarrow \text{man}$ and vice-versa.
 361

362 In the following paragraphs we provide an analysis of the results of our experiments.
 363

364 **Number of queries** One of the biggest limitations of previous methods is that the number of queries
 365 required before a counterexample is found grows exponentially as the hypothesis approximates the
 366 target, as seen in Figure 1. In contrast, the longest running direct EQ run was with the Mistral(24B)
 367 model which required an average of 46.7 queries before termination as seen in Table 3. This is a
 368 dramatic reduction in the number of queries needed.
 369

370 **Consistency of Hypotheses** The standard deviations on the number of queries and in the size of
 371 the intersections show that the direct EQ method has more variability across runs. For the smaller
 372 Mistral(7B) model, terminating after consecutive errors inflates the query count as the model is
 373

374 ⁵A bug in the code was found where the \log_2 was used instead of the natural log. This results in a slightly
 375 lower threshold in the experiments.

376 ⁶www.wikidata.org

377 ⁷During our calculations of the confidence of extracted rules, the Mistral(7B) model using the direct EQ
 approach extracted the rule $\neg(\text{before_1875})$ which we removed before calculating the confidence.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

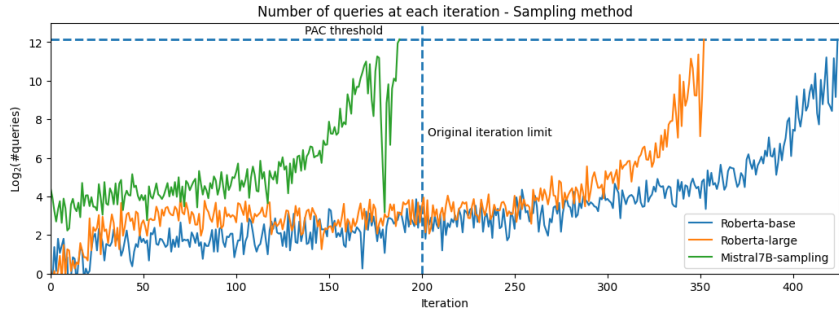


Figure 1: Log_2 of the number of queries per iteration using the sampling method. The number of queries needed grows exponentially as the size of the hypothesis grows. The threshold for termination was set using PAC framework

Table 3: Size of Horn(\mathcal{H}) and non-Horn(\mathcal{Q}) rules, intersections, expected intersections, average weighted confidence, and the number of queries for both the sampling and the direct EQ methods.

| Sampling method | | | | | | |
|------------------|-----------------------|-------------|-------------|-----------|-------------|-----------------------|
| Model | Size(\mathcal{H}) | $ \cap $ | $E(\cap)$ | Conf | #Queries | Size(\mathcal{Q}) |
| RoBERTa-B | 94.00±0.00 | 94.00±0.00 | 9.82±0.00 | 0.48±0.32 | 14,839±1921 | 91.0±0.0 |
| RoBERTa-L | 110.90±0.30 | 110.80±0.40 | 13.67±0.05 | 0.74±0.25 | 13,140±1470 | 63.9±0.3 |
| Mistral(7B) | 55.62±0.78 | 53.85±1.86 | 1.56±0.03 | 0.58±0.46 | 18,271±2245 | 28.5±1.2 |
| Direct EQ method | | | | | | |
| Model | Size(\mathcal{H}) | $ \cap $ | $E(\cap)$ | Conf | #Queries | Size(\mathcal{Q}) |
| Mistral(7B) | 9.70±2.69 | 2.53±1.33 | 0.10±0.04 | 0.50±0.32 | 11.5±10.9 | 0.1±0.4 |
| Mistral(24B) | 15.75±6.66 | 3.42±1.76 | 0.27±0.16 | 0.84±0.34 | 46.7±23.0 | 0.0±0.0 |
| DeepSeek | 3.10±1.67 | 0.15±0.43 | 0.00±0.00 | 0.84±0.34 | 7.8± 9.8 | 0.1±0.2 |

often not able to parse the error description well and the model repeats the same mistake again. This repetitive behaviour was what motivated us to include this error threshold. This premature termination behaviour results in that the size of the extracted hypotheses for the smaller Mistral(7B) model is much smaller when using the direct EQ method. The larger Mistral(24B) model and the Deepseek reasoning model were better able to follow instructions and therefore more likely to terminate intentionally. The Deepseek model often searched the hypothesis space within its <think> tags, rejecting many possible counterexamples in the process. One text extracted from the think tag of a run with the Deepseek model has been provided for illustration purposes and can be found in appendix A.5. This highly selective behaviour likely influenced the small size of the resulting hypotheses and the low consistency for the reasoning model. The size of the intersections across runs are much larger than the expected intersections for both approaches showing that both approaches are able to extract common rules at a rate well above random chance.

Correlation and confidence In Table 4 we see a subset of the rules with the highest confidence extracted with the sampling method, and Table 5 shows the same for models with the direct EQ approach. Both approaches generate very confident rules but the sampling approach also has a tendency to generate many low confidence rules, seen in Tables 6, 7, and 8 in the appendix. This and the sizes of the extracted hypotheses indicates that the sampling approach are more likely to accept an example as a counterexample than the direct EQ approach. This is supported by the size of the non-Horn rules \mathcal{Q} which is generated along with the hypothesis \mathcal{H} by Algorithm 2 for each run. For a rule to appear in \mathcal{Q} , it first needs to be appear in \mathcal{H} , before being proven to be representing a non-Horn clause by later positive counterexamples. From the average size of \mathcal{Q} for the sampling method we see that many rules extracted were later removed from \mathcal{H} by the algorithm. The remaining low confidence rules may be non-Horn rules that were not removed from \mathcal{H} before hitting the PAC threshold. These low confidence rules bring down the weighted average confidence of the extracted rules for the sampling approach.

Table 4: Five most confident rules using the sampling method.

| Model | Rules | Rate | Conf |
|---------------|--|------|------|
| RoBERTa-base | $\neg(\text{Australia} \wedge \text{before 1875} \wedge \text{lawyer} \wedge \text{woman})$ | 1.00 | 1.00 |
| | $\neg(\text{priest} \wedge \text{woman})$ | 1.00 | 0.98 |
| | $\neg(\text{before 1875} \wedge \text{mathematician} \wedge \text{woman})$ | 1.00 | 0.98 |
| | $\neg(\text{between 1875 and 1925} \wedge \text{diplomat} \wedge \text{woman})$ | 1.00 | 0.98 |
| | $\neg(\text{banker} \wedge \text{woman})$ | 1.00 | 0.96 |
| RoBERTa-large | $\neg(\text{before 1875} \wedge \text{diplomat} \wedge \text{woman})$ | 1.00 | 1.00 |
| | $\neg(\text{before 1875} \wedge \text{lawyer} \wedge \text{woman})$ | 1.00 | 0.99 |
| | $\neg(\text{S. America} \wedge \text{before 1875} \wedge \text{woman})$ | 1.00 | 0.99 |
| | $\neg(\text{priest} \wedge \text{woman})$ | 1.00 | 0.98 |
| | $\neg(\text{before 1875} \wedge \text{mathematician} \wedge \text{woman})$ | 1.00 | 0.98 |
| Mistral(7B) | $\neg(\text{Oceania} \wedge \text{before 1875} \wedge \text{mathematician} \wedge \text{woman})$ | 1.00 | 1.00 |
| | $\neg(\text{Australia} \wedge \text{before 1875} \wedge \text{man} \wedge \text{nurse})$ | 1.00 | 1.00 |
| | $(\text{Asia} \wedge \text{before 1875} \wedge \text{footballer}) \rightarrow \text{man}$ | 1.00 | 1.00 |
| | $(\text{Australia} \wedge \text{before 1875} \wedge \text{diplomat}) \rightarrow \text{man}$ | 1.00 | 1.00 |
| | $\neg(\text{Africa} \wedge \text{between 1875 and 1925} \wedge \text{footballer} \wedge \text{woman})$ | 1.00 | 1.00 |

Table 5: Five most confident rules using the direct EQ method.

| Model | Rule | Rate | Conf |
|--------------|---|------|------|
| Mistral(7B) | $\neg(\text{N. America} \wedge \text{priest} \wedge \text{woman})$ | 0.05 | 0.96 |
| | $\neg(\text{Oceania} \wedge \text{man} \wedge \text{nurse})$ | 0.05 | 0.96 |
| | $\neg(\text{S. America} \wedge \text{between 1925 and 1951} \wedge \text{mathematician} \wedge \text{woman})$ | 0.05 | 0.87 |
| | $\neg(\text{Africa} \wedge \text{diplomat} \wedge \text{woman})$ | 0.05 | 0.85 |
| | $\neg(\text{Europe} \wedge \text{after 1970} \wedge \text{dancer} \wedge \text{man})$ | 0.05 | 0.59 |
| Mistral(24B) | $\neg(\text{Africa} \wedge \text{before 1875} \wedge \text{man} \wedge \text{nurse})$ | 0.85 | 1.00 |
| | $\neg(\text{S. America} \wedge \text{banker} \wedge \text{before 1875} \wedge \text{woman})$ | 0.25 | 1.00 |
| | $\neg(\text{Africa} \wedge \text{banker} \wedge \text{before 1875} \wedge \text{woman})$ | 0.20 | 1.00 |
| | $\neg(\text{S. America} \wedge \text{before 1875} \wedge \text{mathematician} \wedge \text{woman})$ | 0.15 | 1.00 |
| | $\neg(\text{Australia} \wedge \text{before 1875} \wedge \text{footballer} \wedge \text{woman})$ | 0.10 | 1.00 |
| Deepseek(8B) | $\neg(\text{Europe} \wedge \text{before 1875} \wedge \text{priest} \wedge \text{woman})$ | 0.11 | 1.00 |
| | $\neg(\text{Americas} \wedge \text{between 1875 and 1925} \wedge \text{footballer} \wedge \text{woman})$ | 0.05 | 1.00 |
| | $\neg(\text{Eurasia} \wedge \text{before 1875} \wedge \text{lawyer} \wedge \text{woman})$ | 0.05 | 1.00 |
| | $\neg(\text{Europe} \wedge \text{before 1875} \wedge \text{lawyer} \wedge \text{woman})$ | 0.11 | 1.00 |
| | $\neg(\text{Europe} \wedge \text{between 1875 and 1925} \wedge \text{priest} \wedge \text{woman})$ | 0.11 | 0.99 |

7 CONCLUSIONS

In this paper we explore a new strategy for extracting Horn rules from large language models. We propose an evaluation scheme for the extracted Horn rules using three criteria and compare results with previously proposed strategies used on the BERT family of models. We find that our new strategy significantly reduces the number of queries required while extracting on average higher confidence rules. The significant number of queries required with the sampling method makes it infeasible for use with generative LLMs. Our new strategy thus provides an alternative method for extracting Horn rules from generative LLMs.

Limitations and Future Work A limitation of this work is the scope of the experiments; the number and type of models tested. While the case study was meant to showcase an application of the strategy, we encountered challenges with prompt sensitivity, leading to diverging results, as well as challenges interpreting responses as smaller models were not able to process the instructions in the prompt and provide valid counterexamples as the set of extracted rules grew larger. Another limitation of this work is that we do not address how the non-deterministic behaviour of generative models impact the guarantees of Algorithm 2. As future work, we would like to investigate other case studies since the approach is independent of the variables in Table 2, improve the expressivity of the rule language, and provide a formal treatment for the changes in the answers of the LLMs.

486 7.1 ETHICS STATEMENT
487

488 This work extracts rules on gender-occupation relationships in LLMs. We compare the rules extracted
489 against historical data which in some cases reflect historical gender biases. We wish to acknowledge
490 that when working with such biases it is important to emphasize that it is not our intention to
491 perpetuate historical biases. With this work, we only wish to contribute to the removal of harmful
492 biases by highlighting their existence.

493 7.2 REPRODUCIBILITY STATEMENT
494

495 The code used to run the experiments is added as supplementary material. A description of each
496 file in the repository is included in the README.md. All experiments were performed on a
497 High-Performance Computing (HPC) cluster with 1 GPU-accelerated node per experiment. The
498 specific GPUs were either a NVIDIA A100 or a NVIDIA RTX3090, limited to 90GB of memory
499 per run. The historical data used to compare the results were extracted from, and is available on
500 www.wikidata.org. The code for the extraction of the data is available in the code in the
501 supplementary material.
502

503 REFERENCES
504

- 505 Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of
506 items in large databases. In *Proceedings of the 1993 ACM SIGMOD International Conference*
507 *on Management of Data*, SIGMOD '93, pp. 207–216, New York, NY, USA, 1993. Association
508 for Computing Machinery. ISBN 0897915925. doi: 10.1145/170035.170072. URL <https://doi.org/10.1145/170035.170072>.
509
- 510 Dana Angluin. Queries and concept learning. *Mach. Learn.*, 2(4):319–342, 1987. doi: 10.1007/
511 BF00116828. URL <https://doi.org/10.1007/BF00116828>.
512
- 513 Dana Angluin, Michael Frazier, and Leonard Pitt. Learning conjunctions of horn clauses. *Mach.*
514 *Learn.*, 9:147–164, 1992. doi: 10.1007/BF00992675. URL [https://doi.org/10.1007/
515 BF00992675](https://doi.org/10.1007/BF00992675).
516
- 517 Sophie Blum, Raoul Koudijs, Ana Ozaki, and Samia Touileb. Learning horn envelopes via queries
518 from language models. *Int. J. Approx. Reason.*, 171:109026, 2024. doi: 10.1016/J.IJAR.2023.
519 109026. URL <https://doi.org/10.1016/j.ijar.2023.109026>.
520
- 521 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
522 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
523 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
524 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
525 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
526 Sutskever, and Dario Amodei. Language Models are Few-Shot Learners, July 2020. URL
<http://arxiv.org/abs/2005.14165>.
- 527 Xinyun Chen, Ryan A. Chi, Xuezhi Wang, and Denny Zhou. Premise Order Matters in Reasoning
528 with Large Language Models, May 2024. URL <http://arxiv.org/abs/2402.08939>.
529 arXiv:2402.08939 [cs].
530
- 531 Rina Dechter and Judea Pearl. Structure identification in relational data. *Artificial Intelli-*
532 *gence*, 58(1):237–270, 1992. ISSN 0004-3702. doi: [https://doi.org/10.1016/0004-3702\(92\)](https://doi.org/10.1016/0004-3702(92)90009-M)
533 [90009-M](https://www.sciencedirect.com/science/article/pii/000437029290009M). URL [https://www.sciencedirect.com/science/article/pii/
534 000437029290009M](https://www.sciencedirect.com/science/article/pii/000437029290009M).
- 535 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,
536 2025. URL <https://arxiv.org/abs/2501.12948>.
537
- 538 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep
539 bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL
<http://arxiv.org/abs/1810.04805>.

- 540 Jean-Louis Guigues and Vincent Duquenne. Familles minimales d’implications informatives résultant
541 d’un tableau de données binaires. *Mathématiques et Sciences humaines*, 95:5–18, 1986.
- 542
- 543 Qiyuan He, Yizhong Wang, Jianfei Yu, and Wenya Wang. Language models over large-scale
544 knowledge base: on capacity, flexibility and reasoning for new facts. In Owen Rambow, Leo
545 Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert
546 (eds.), *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 1736–
547 1753, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL <https://aclanthology.org/2025.coling-main.118/>.
- 548
- 549 Jie Huang and Kevin Chen-Chuan Chang. Towards Reasoning in Large Language Models: A Survey,
550 May 2023. URL <http://arxiv.org/abs/2212.10403>. arXiv:2212.10403 [cs].
- 551
- 552 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
553 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
554 Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
555 Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- 556
- 557 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
558 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining
559 approach. *CoRR*, abs/1907.11692, 2019. URL <http://arxiv.org/abs/1907.11692>.
- 560
- 561 Philipp Mondorf and Barbara Plank. Beyond Accuracy: Evaluating the Reasoning Behavior of Large
562 Language Models – A Survey, August 2024. URL <http://arxiv.org/abs/2404.01869>.
arXiv:2404.01869 [cs].
- 563
- 564 Ana Ozaki, Roberto Confalonieri, Ricardo Guimarães, and Anders Imenes. Extracting PAC decision
565 trees from black box binary classifiers: The gender bias study case on bert-based language models.
566 In Toby Walsh, Julie Shah, and Zico Kolter (eds.), *AAAI*, pp. 19767–19775. AAAI Press, 2025.
567 doi: 10.1609/AAAI.V39I18.34177. URL <https://doi.org/10.1609/aaai.v39i18.34177>.
- 568
- 569 Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu,
570 and Alexander Miller. Language models as knowledge bases? In Kentaro Inui, Jing Jiang,
571 Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods
572 in Natural Language Processing and the 9th International Joint Conference on Natural Language
573 Processing (EMNLP-IJCNLP)*, pp. 2463–2473, Hong Kong, China, November 2019. Association
574 for Computational Linguistics. doi: 10.18653/v1/D19-1250. URL <https://aclanthology.org/D19-1250/>.
- 575
- 576 Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the
577 parameters of a language model? In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu
578 (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing
579 (EMNLP)*, pp. 5418–5426, Online, November 2020. Association for Computational Linguistics.
580 doi: 10.18653/v1/2020.emnlp-main.437. URL <https://aclanthology.org/2020.emnlp-main.437/>.
- 581
- 582 Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu,
583 Mingyu Ding, Hongyang Li, Mengzhe Geng, Yue Wu, Wenhai Wang, Junsong Chen, Zhangyue
584 Yin, Xiaozhe Ren, Jie Fu, Junxian He, Wu Yuan, Qi Liu, Xihui Liu, Yu Li, Hao Dong, Yu Cheng,
585 Ming Zhang, Pheng Ann Heng, Jifeng Dai, Ping Luo, Jingdong Wang, Ji-Rong Wen, Xipeng Qiu,
586 Yike Guo, Hui Xiong, Qun Liu, and Zhenguo Li. A survey of reasoning with foundation models,
587 2024. URL <https://arxiv.org/abs/2312.11562>.
- 588
- 589 Leslie G. Valiant. A theory of the learnable. In Richard A. DeMillo (ed.), *Proceedings of the 16th
590 Annual ACM Symposium on Theory of Computing, April 30 - May 2, 1984, Washington, DC,
591 USA*, pp. 436–445. ACM, 1984. doi: 10.1145/800057.808710. URL <https://doi.org/10.1145/800057.808710>.
- 592
- 593 Gail Weiss, Yoav Goldberg, and Eran Yahav. Extracting automata from recurrent neural networks
using queries and counterexamples (extended version). *Mach. Learn.*, 113(5):2877–2919, 2024.
doi: 10.1007/S10994-022-06163-2.

594
595
596
597
598
599
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

Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. Satlm: Satisfiability-aided language models using declarative prompting, 2023. URL <https://arxiv.org/abs/2305.09656>.