Can We Infer Confidential Properties of Training Data from LLMs?

Pengrun Huang, Chhavi Yadav, Kamalika Chaudhuri[†], Ruihan Wu[†]

University of California, San Diego {peh006, cyadav, kamalika, ruw076}@ucsd.edu

Abstract

Large language models (LLMs) are increasingly fine-tuned on domain-specific datasets to support applications in fields such as healthcare, finance, and law. These fine-tuning datasets often have sensitive and confidential dataset-level properties — such as patient demographics or disease prevalence—that are not intended to be revealed. While prior work has studied property inference attacks on discriminative models (e.g., image classification models) and generative models (e.g., GANs for image data), it remains unclear if such attacks transfer to LLMs. In this work, we introduce PropInfer, a benchmark task for evaluating property inference in LLMs under two fine-tuning paradigms: question-answering and chat-completion. Built on the ChatDoctor dataset, our benchmark includes a range of property types and task configurations. We further propose two tailored attacks: a prompt-based generation attack and a shadow-model attack leveraging word frequency signals. Empirical evaluations across multiple pretrained LLMs show the success of our attacks, revealing a previously unrecognized vulnerability in LLMs. We release our code at github.com/PengrunH/Property_inference_attack_LLM.

1 Introduction

Large language models (LLMs) are increasingly deployed in real-world applications across domains such as healthcare [14], finance [19], and law [17]. To adapt to domain-specific tasks, such as customer service or tele-medicine, these models are typically fine-tuned on proprietary datasets that are relevant to the tasks at hand before deployment. These domain-specific fine-tuning datasets however often contain *dataset-level confidential information*. For example, a customer-service dataset sourced from a business may contain information about their typical customer-profile; a doctor-patient chat dataset sourced from a hospital may contain patient demographics or the fraction of patients with a sensitive disease such as HIV. Many businesses and medical practices would consider this kind of information non-public for business or other reasons. Thus, unintentional leakage of this information through a deployed model could lead to a breach of confidentiality. Unlike individual-level privacy breaches that is typically addressed by rigorous definitions such as differential privacy [9, 10], the risk here is the leakage of dataset-level properties.

Prior work has investigated this form of leakage, commonly referred to as *property inference* [2]. Most of the literature here has focused on two settings. The first involves discriminative models trained on tabular or image data [2, 11, 6, 28, 33, 12], where the goal is to infer attributes such as the gender distribution in a hospital dataset. The second focuses on generative models [34, 31], such as GANs for face synthesis, where attackers may attempt to recover aggregate properties such as the racial composition of the training data. In both cases, property inference has been shown to be feasible, and specialized attacks have been proposed to exploit these vulnerabilities.

^{1†} Equal advising.

However, property inference in large language models (LLMs) introduces two distinct challenges. First, unlike inferring a single attribute from models trained on tabular data, the sensitive properties are more complex and may be indirectly embedded within the text. For example, gender might be implied through broader linguistic cues, such as the mention of a "my gynecologist". LLMs may memorize such properties implicitly, making them more challenging to infer reliably in property inference studies. The second challenge is that, unlike the models typically studied in prior work, LLMs do not fit cleanly into purely discriminative or generative categories; this raises questions about what kind of property inference attacks apply and succeed for these problems.

In this work, we investigate both questions by introducing a new benchmark task – PropInfer¹ – for property inference in LLMs. Our task is based on the Chat-Doctor dataset [18] – a domain-specific medical dataset containing a collection of question-answer pairs between patients and doctors. There are two standard ways to fine-tune an LLM with this dataset that correspond to two use-cases: question-answering and chat-completion. According to the use case, our benchmark task has two modes where the models are fine-tuned differently – Q&A Mode and Chat-Completion Mode. To comprehensively study property inference across the two modes of models, we select a range of properties that are explicitly or implicitly reflected in both questions and answers.

We propose two property inference attacks tailored to LLMs. The first is a black-box generation-based attack, inspired by prior work [34]; the intuition is that the distribution of the generated samples reflect the distribution of the fine-tuning data. Given designed prompts that reflects characteristics of the target dataset, the adversary generates multiple samples from the target LLM and labels each based on the presence of the target property. The property ratio is then estimated by aggregating the labels. The second is a shadow-model attack with word-frequency. With access to an auxiliary dataset, the adversary first trains a set of shadow models with varying property ratios and extracts word frequency from the shadow models based on some selected keyword list. Then the adversary trains a meta-attack model that maps these frequencies to the corresponding property ratios. This enables the inference on the target model by computing its output word frequencies.

We empirically evaluate our two attacks alongside baseline methods using our PropInfer-benchmark. Our results show that the shadow-model attack with word frequency is particularly effective when the target model is fine-tuned in the Q&A Mode and the target property is more explicitly revealed in the question content than the answer. In contrast, when the model is fine-tuned in Chat-Completion Mode or when the target attribute are embedded in both the question and the answer, the black-box generation-based attack proves to be simple yet highly effective.

Our experimental results reveal a previously underexplored vulnerability in large language models: property inference, which enables adversaries to extract dataset-level attributes from fine-tuned models. This finding exposes a tangible threat to data confidentiality in real-world deployments. It also underscores the need for robust defense mechanisms to mitigate such attacks – an area where our benchmark provides a standardized and extensible framework for future research and evaluation.

2 Related Work

Property inference. Property Inference Attack (PIA) was first described by [2], as follows: given two candidate training data distributions $\mathcal{D}_1, \mathcal{D}_2$ and a target model, the adversary tries to guess which training distribution (out of $\mathcal{D}_1, \mathcal{D}_2$) is the target model trained on. Typically, the two candidate distributions only differ in the marginal distribution of a binary variable, such as gender ratio. A major portion of past work on property inference focuses on discriminative models [2, 11, 6, 28, 33, 12]; here the attacks mainly rely on training meta-classifiers on some representations to predict target ratio. For example, in the white-box setting, [2, 11] use model weights as the input of the meta-classifier to predict the correct distribution. In the grey-box setting, where the adversary have access to the training process and some auxiliary data, [27, 28] use model outputs such as loss or probability vector as inputs to the meta-classifier.

Moving on to generative models, [34] study property inference attack for GANs. The target GANs are trained on a human-face image dataset, whereas the adversary's task is to predict the ratio of the target property among the dataset, such as gender or race. Their attack follows the intuition that the

¹https://huggingface.co/datasets/Pengrun/PropInfer_dataset

generated samples from GANs can reflect the training distribution. Later on, [31] studies property existence attacks. For example, if any images of a specific brand of cars are used in the training set.

Contrary to previous works which either focus on discriminative models or pure generative models, we consider property inference attack for large language models. Since the model architecture, training paradigm and data type for LLMs are very distinct from previous works, it is unclear whether previous attacks still apply in the LLM setting.

Other related works on data privacy and confidentiality in LLMs. [5] study training data extraction from LLMs, aiming to recover individual training samples. While one might try to infer dataset properties from extracted data, this often fails because the extracted samples are typically biased and not representative of the overall distribution. [20] investigate dataset inference attacks, which aim to identify the dataset used for fine-tuning from a set of candidates. In contrast, our goal is to infer specific aggregate properties, not the dataset itself. [26] study idiosyncrasies in public LLMs, determining which public LLM is behind a black-box interface. Although they also use word frequency signals, their objective differs fundamentally from ours.

3 Preliminaries

3.1 Large Language Models Fine-Tuning

A large language model (LLM) predicts the likelihood of a sequence of tokens. Given input tokens $t_0,...,t_{i-1}$, the language model parameterized by parameters $\theta,\,f_\theta,\,$ outputs the distribution of the possible next token $f_\theta(t_i|t_0,...,t_{i-1})$. Pre-training LLMs on large-scale corpora enables them to develop general language understanding and encode broad world knowledge. In pre-training, the LLM is trained to maximize the likelihood of unlabeled text sequences. Each training sample is a document comprising a sequence of tokens, and the objective is to minimize the negative log-likelihood: $\mathcal{L}(\theta) = -\sum_{i=1}^k \log f_\theta(t_i|t_0,...t_{i-1})$ where k is the total number of tokens in the document.

After pre-training, LLMs are often fine-tuned on domain-specific datasets to improve performance on downstream tasks. The data in such datasets typically consists of an instruction (I), which generally describes the task in natural language, and a pair of an input (x) and a ground-truth output (y). Two popular fine-tuning approaches are:

- Supervised Fine-Tuning (SFT; [24, 22]). SFT minimizes the negative log-likelihood of the output tokens conditioned on the instruction and input. This approach focuses on learning the mapping from (I, x) to y and is commonly used in question-answering tasks. The objective is:
 \$\mathcal{L}_{\text{SFT}}(\theta) = -\sum_{i=1}^{l} \log f_\theta(y_i | I, x, y_0, ..., y_{i-1}).\$

 Causal Language-Modeling Fine-Tuning (CLM-FT; [24]). Different from SFT, CLM-FT
- 2. Causal Language-Modeling Fine-Tuning (CLM-FT; [24]). Different from SFT, CLM-FT follows the pre-training paradigm and minimizes the loss over all tokens in the concatenated sequence t of instruction, input, and output (I,x,y). This method treats the full sequence autoregressively, making it suitable for tasks involved auto-completion for both user and chatbot. The objective is $\mathcal{L}_{\text{CLM-FT}}(\theta) = -\sum_{i=1}^k \log f_{\theta}(t_i|t_0,...t_{i-1})$.

3.2 Property Inference Attack

In this paper, we focus on LLMs that have been fine-tuned on domain-specific datasets, as these datasets often encompass scenarios involving confidential or sensitive information. Given an LLM fine-tuned on such a dataset, property inference attacks aim to extract the dataset-level properties of the fine-tuning dataset from the finetuned LLM, which the data owner does not intend to disclose ².

Let $\mathcal{S}=(x_i,y_i)_{i=1}^n$ denote the fine-tuning dataset of size n, consisting of i.i.d. samples drawn from an underlying distribution \mathcal{D} over the domain $X\times Y$. We denote the fine-tuned model as $f=\mathcal{A}(\mathcal{S};I)$, where \mathcal{A} is the fine-tuning algorithm applied to \mathcal{S} using a fixed instruction template I. Let $P:X\times Y\to \{0,1\}$ be a binary function indicating whether a particular data point satisfies a certain property. For example, P(x,y)=1 may indicate that a patient in a doctor-patient dialogue (x,y) is female. The adversary's goal is to estimate the ratio of the target property P among the dataset S. The adversary's goal is: $r(P,S):=\frac{1}{n}\sum_{i=1}^n P(x_i,y_i)$.

²When the property is correlated with what the model learns, it seems pessimistic to avoid such leakage. However, the properties of concern in practice are often orthogonal to the task itself. See details in Section 7.

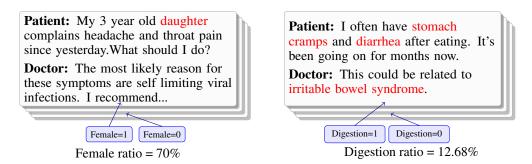


Figure 1: This figure demonstrates examples of the ChatDoctor dataset and the property labels. (Left) An example of dialogue explicitly indicate the patient is a female, since it mentioned "daughter"; (right) an example of dialogue indicating the patient is consulting about digestive disorder.

3.3 Threat Models

We consider two standard threat models in this work: black-box setting and grey-box setting.

Black-box setting. Following prior work [34], we consider the black-box setting in which the adversary has only API-level access to the target model f. In the LLM context, this means the adversary can create arbitrary prompts and receive sampled outputs from the model, but has no access to its parameters, architecture, or the auxiliary data. This represents the most restrictive and least informed setting for the adversary, where only input-output interactions are observable.

Grey-box setting. Another standard threat model in the literature is this grey-box access [34, 27, 28]. In addition to black-box access to the target model f, we assume the adversary (1) has knowledge of the fine-tuning procedure \mathcal{A} , including details of the pre-trained model, fine-tuning method and the instruction template I, (2) has the knowledge of target dataset size n, and (3) has an auxiliary dataset $\mathcal{S}_{\text{aux}} = (\hat{x}_i, \hat{y}_i)_{i=1}^{n'}$ drawn i.i.d. from the underlying distribution $\hat{\mathcal{D}}$. The inference problem becomes trivial when $\hat{\mathcal{D}}$ is the same as \mathcal{D} where the fine-tuning dataset \mathcal{S} is sampled from. To make the setting nontrivial and realistic, we assume that $\hat{\mathcal{D}}$ and \mathcal{D} differ only in the marginal distribution of the target property, while sharing the same conditional distribution given the property.

4 PropInfer: Benchmarking Property Inference Across Fine-Tuning and Property Types

We build our benchmarks upon a popular patient-doctor dialogues dataset ChatDoctor [18]; Figure 1 shows examples of the dataset. In this setting, an adversary may attempt to infer sensitive demographic attributes or the frequency of specific medical diagnoses — both representing realistic threats in which leakage of aggregate properties could have serious consequences. To systematically study property inference attacks in LLMs, we extend the original ChatDoctor dataset by introducing two modes of the fine-tuned models, and the target properties, into our benchmark.

Two modes of the fine-tuned models: Q&A Mode and Chat-Completion Mode. The ChatDoctor dataset supports two common use cases: (1) Doctor-like Q&A chatbot for automatic diagnosis, and (2) Chat-completion to assist both patients and doctors. Let x denote the patient's symptom description and y the doctor's diagnosis. In the Q&A chatbot mode, models are fine-tuned using Supervised Fine-Tuning (SFT), learning to generate y conditioned on I and x. In the chat-completion mode, models are trained using Causal Language-Modeling fine-tuning (CLM-FT), which minimizes loss over the entire sequence of tokens in I, x, and y. This allows the model to predict tokens at any point in the dialogue. For the formal training objectives, kindly refer to Section 3.1. Accordingly, our benchmark includes both the Q&A Mode and the Chat-Completion Mode, reflecting two widely used fine-tuning paradigms: SFT and CLM-FT.

These two modes naturally introduce different memorization patterns: CLM-FT encourages the model to learn the joint distribution $\mathbb{P}(I, x, y)$, potentially memorizing both patient and doctor texts equally; differently, SFT focuses on the conditional distribution $\mathbb{P}(y|I, x)$, emphasizing the doctor's

response y more heavily than the patient's input x. Consequently, effective attack strategies may differ across two fine-tuning modes, motivating separate analyses in our benchmark.

Target properties: the demographic information and the medical diagnosis frequency. Since two fine-tuned modes have different memorization patterns, property inference behavior can vary depending on where the target property resides. We therefore propose two categories of the properties: the demographic information, which is often revealed in the patient description, and the medical diagnoses, which are discussed by both the patients and the doctor, as shown in Figure 1.

For demographic information, we select patient *gender*, which can be explicitly stated (e.g., "I am female") or implicitly suggested (e.g., "pregnancy" or "periods") in patient descriptions x. We label the gender property using ChatGPT-40 and filter out samples with ambiguous gender indications. This results in a gender-labeled dataset of 29,791 conversations, in which 19,206 samples have female labels and 10,585 have male labels. We use 15,000 samples to train the target models and the remaining 14,791 as auxiliary data for evaluating attacks in the grey-box setting. For medical diagnosis attributes, we use the original training split of the ChatDoctor dataset with size 50,000 for training the target models and consider three binary properties: (1) Mental disorders (5.10%), (2) Digestive disorders (12.68%), and (3) Childbirth (10.6%). Please see Appendix A.1 for details on the labeling process and Section 6.1 for details on task definitions and model fine-tuning procedures.

5 Attacks

Recall that the goal of the adversary is to estimate the value of the property of interest for the target model $M_{\rm target}$. Prior work [34, 27] has given attacks that can achieve this goal on simpler models or image and tabular data, and therefore these do not apply directly to the LLM setting. Inspired by the initial ideas from the old attacks, we proposes two new attacks tailored for the LLM setting.

5.1 Generation-Based Attack under Black-Box Setting

Prior work [34] introduced an output-generation-based property inference attack under black-box access, specifically targeting unconditional GANs. In the context of LLMs, which perform conditional token-level generation, we adapt this approach by generating outputs based on carefully designed input prompts that constrain the generation distribution to the domain of interest. Our adapted attack consists of the following three steps.

Prompt-conditioned generations. We construct a list of prompts T, that encodes high-level contextual information about the fine-tuning dataset. For example, for the Chat-Doctor dataset, a prompt like "Hi, doctor, I have a medical question." would be a reasonable choice. Given any prompt $t \in T$, we generate a corresponding set of output samples $S_{f,t}$ from the target model f.

Property labeling. We define a property function \hat{P} hold by the adversary, which maps each generated sample $s \in S_{f,t}$ to a value in $\{0,1,N/A\}$. A label of 1 or 0 indicates whether the sample reflects the presence or absence of the target property respectively. \hat{P} assigns the label N/A for the samples that are ambiguous or indeterminate with respect to the property.

Prompt-based property inference. To estimate the property ratio, we first restrict attention to samples with valid labels. Let $S_{f,t}^* \subseteq S_{f,t}$ denote the subset of generated samples for which $\hat{P}(s) \neq N/A$. The estimated ratio given the prompt t is $\hat{r}_t = \frac{1}{|S_{f,t}^*|} \sum_{s \in S_{f,t}^*} \hat{P}(s)$. If the adversary uses a list of prompts T, the aggregated estimation across prompts is given by: $\hat{r} = \frac{1}{|T|} \sum_{t \in T} \hat{r}_t$.

5.2 Shadow-Model Attack with Word Frequency under Grey-Box Setting

Prior work [27, 28, 13] has proposed various shadow-model based property inference attacks. The core idea is that the adversary trains multiple shadow models on an auxiliary dataset that is disjoint from the target model's dataset, with varying target property ratios. Given both the shadow models and their ground-truth property ratios, the adversary can learn a mapping from some extracted model features to the underlying property ratios. The framework ³ is describes as follows:

³Prior work frames property inference as a hypothesis testing problem between two candidate ratios. Our framework extends the existing framework by enabling the adversary to predict property ratios directly.

- 1. **Shadow model training.** The adversary selects k_1 target property ratios $r_1, \ldots, r_{k_1} \in [0, 1]$. For each ratio r_i , the adversary subsamples k_2 auxiliary datasets to match r_i with the target size n, and fine-tunes LLMs with the same fine-tuning procedure \mathcal{A} , resulting in $k_1 \cdot k_2$ shadow models. The shadow models can be denoted as $f_{i,j}$, where i indexes the ratio, and j indexes the repetition.
- 2. Meta attack model training through a defined shodow feature function. A shadow feature function F maps each model to a d-dimensional feature vector. Given the shadow models and their corresponding ratios, a meta dataset is constructed: $(F(f_{i,j}), r_i) \mid i \in [k_1], j \in [k_2]$. A meta attack model $g: \mathbb{R}^d \to [0,1]$ is learned from the meta dataset to predict the property ratio from the extracted model features. In this paper, we use XGBoost [7] to train the meta attack model.
- 3. **Property inference.** The final inference on the target model f is made by computing $\hat{r} = g(F(f))$.

Constructing new shadow attacks with word frequency. The choice of the shadow feature functions F plays an important role in the success of the attack. While previous work relies on loss or probability vector [27, 28], some studies have shown that these features may not be the most effective way to measure the performance of the LLMs [5, 8]. We alternatively focus on another feature specific to the LLM setting, word frequency, which has been parallelly studied for the literature of membership inference attack [16, 21]. Our attack is based on the intuition that certain properties may strongly correlate with the appearance of specific words in the text. As a result, models fine-tuned on datasets with different property distributions may exhibit distinct word patterns in their generations.

Assume V^* is a selected list of d keywords, which we will describe its construction later. Similar to the generation attack, given a model f and the prompt t that describes the meta information about the fine-tuning dataset, we generate a set of text samples $S_{f,t}$. For each word $v \in V^*$, we calculate the word-frequency $\mu_v^{f,t}$, defined as the proportion of samples in $S_{f,t}$ containing v. If the adversary uses a list of prompts, it can average this by $u_v^f = \frac{1}{T} \sum_{t \in T} u_v^{f,t}$. The resulting vector $(\mu_v^f)_{v \in V^*} \in [0,1]^d$ serves as the shadow feature, and the shadow feature function is defined as $F_{\text{word}}(f) := (\mu_v^f)_{v \in V^*}$.

To construct the keyword list V^* , we first define the full vocabulary V as all words that appear in at least one sample in any $S_{f_{i,j},t}$. Then we apply a standard feature selection algorithm ⁴ using the word frequency $(\mu_v^{f_{i,j}})_{v\in V}$ and their corresponding labels(i.e. the property ratios). This process selects the d most informative words for the property ratio prediction task, forming the final keyword list V^* .

6 Experiments

In this section, we empirically evaluate the effectiveness of our proposed attacks within the newly introduced benchmark, PropInfer. Specifically, we aim to answer the following research questions:

- 1. How do the proposed attacks perform in Chat-Completion Mode versus Q&A Mode?
- 2. How does the choice of fine-tuning method influence the success of property inference attacks?

6.1 Experimental Setup

For implementation details, including the selection of hyperparameters for fine-tuning, our attacks, and baseline methods, please refer to Appendix A.2.

Models. We use three open base models for experimentation: Llama-1-8b[30], Pythia-v0-6.9b[3] and Llama-3-8b-instruct [1]. We use the Llama-1 and Pythia-v0 since these were released before the original ChatDoctor dataset and hence have no data-contamination from the pre-training stage, giving us a plausibly more reliable attack performance. While Llama-3 came after ChatDoctor release, we still use it since it is highly performant and is widely used for experimentation. Refer to Appendix A.2 for implementation details and fine-tuning performance.

Property inference tasks. Our benchmark defines two property inference tasks. **Gender property inference**, where the goal is to infer the ratios of female samples in the fine-tuning dataset. We define 3 target ratios of female: $\{0.3, 0.5, 0.7\}$; for each target ratio, we subsample 3 datasets with different random seeds to match each target ratio while keeping the same size 6500, and we evaluate this by attacking the total 9 target models. **Medical diagnosis property inference**, where the goal is to infer the proportion of three diagnosis-related properties (e.g., mental disorder (5.10%), digestive

⁴We used the algorithm f_regression implemented in scikit-learn library [23].

Table 1: Attack Performance for gender property in the Q&A mode and Chat-Completion mode. Reported numbers are the Mean Absolute Errors (MAE; \downarrow) between the predicted and target ratios. We highlight the attack that achieves the smallest total MAE across different target ratios.

| Model | Attacks | Q&A Mode | | | Chat-Completion Mode | | |
|-----------|----------------------------|------------------|-----------------------------|--------------------|----------------------|-----------------------------|-----------------------------|
| | | 30 | 50 | 70 | 30 | 50 | 70 |
| | Direct asking | 23.17±1.78 | $3.98{\pm}1.88$ | 18.6±0 | 22.8±1.98 | 7.7±1.8 | 18.57±4.71 |
| Llama-1 | BB generation | 36.52 ± 0.11 | 15.45 ± 3.09 | 1.45 ± 0.64 | 1.73 ± 0.76 | 2.64 ± 3.33 | 3.28 ± 3.64 |
| Liailia-i | Perplexity | 28.67 ± 9.34 | 9.38 ± 8.95 | 24.16 ± 2.45 | 35.19 ± 10.99 | 14.5 ± 5.98 | $5.33{\pm}6.09$ |
| | Word-frequency | 11.43±3.0 | $7.33{\scriptstyle\pm6.59}$ | 6.85 ± 5.03 | 3.44 ± 4.61 | $0\pm o$ | $6.6{\scriptstyle\pm9.35}$ |
| | Direct asking ⁵ | _ | _ | _ | _ | _ | _ |
| D-41-10 | BB generation | 46.75 ± 3.64 | 23.45 ± 5.89 | 10.31 ± 4.85 | 3.56 ± 2.03 | 5.61 ± 0.78 | 2.15 ± 2.45 |
| Pythia-v0 | Perplexity | $22.33{\pm}15.8$ | 11.25 ± 13.68 | 25.79 ± 15.59 | 4.32 ± 3.25 | 9.94 ± 0.59 | 9.39 ± 0 |
| | Word-frequency | 22 ± 10.4 | $7.95{\scriptstyle\pm9.44}$ | 9.25 ± 11.9 | 3.31 ± 4.68 | $3.27{\scriptstyle\pm4.62}$ | $6.73{\pm}8.22$ |
| Llama-3 | Direct asking | 14.27±5.32 | 4.86±1.33 | 19.9±4.24 | 17.97±5.33 | 4.0 ± 2.12 | 16.17±1.18 |
| | BB generation | 23.64 ± 5.82 | 5.79 ± 6.46 | 14.01 ± 1.68 | 0.61 ± 0.77 | $1.33{\pm}1.31$ | $1.25{\pm}1.52$ |
| | Perplexity | 13.28 ± 4.77 | $25.0{\scriptstyle\pm25.4}$ | 19.01 ± 20.52 | 17.80 ± 9.06 | 19.85 ± 7.6 | 6.24 ± 7.57 |
| | Word-frequency | $8.29{\pm}2.13$ | $7.33{\pm}6.59$ | $10.66 {\pm} 7.12$ | $2.45{\pm}2.3$ | $3.33{\pm}4.7$ | $5.83{\scriptstyle\pm1.73}$ |

disorder (12.68%), childbirth (10.6%) from the medical diagnosis dataset (with size 50,000). We train 3 target models on the entire dataset for evaluation.

For both tasks, we evaluate the attacks on Q& A Mode and Chat-Completion Mode. For the gender inference task, we evaluate both black-box and grey-box attacks, where our benchmark provides auxilary dataset of size 14,791. For the medical diagnosis task, we evaluate only the black-box adversary, as the grey-box setting requires that the auxiliary dataset shares the same conditional distribution given the target property while differing only in the marginal distribution. Constructing a well-matched auxiliary dataset for multiple properties simultaneously is inherently nontrivial.

Our attack setups. For the black-box generation-based attack (BB generation) as described in Section 5.1 on our benchmark, one example of the prompts we used is to fill out the sentence: "Hi, Chatdoctor, I have a medical question." In total, we use three prompts; the full list is provided in Appendix A.2. For each target model f and prompt t, we generate 2000 samples. Each generated text is then labeled by ChatGPT-40 (\hat{P}) based on the target property.

For the **shadow-model attack with word frequency (word-frequency attack)**, as described in Section 5.2, we choose $k_1=7$ property ratios in $\{0.2,0.3,\cdots,0.8\}$, with $k_2=5$ or 6 (varying between different LLMs) shadow models trained per ratio. We apply the same three prompts as in the BB generation and generate ~ 100 k samples for each prompt to estimate the word frequency.

Baseline attacks. We consider three baseline attacks and put some implementation details in Appendix A.2. (1) **Direct asking** (black-box baseline) is a direct query approach, where the adversary simply asks the model to report the property ratio. For example, we prompt the model with: "what is the percentage of patient having mental disorder concern in the ChatDoctor dataset?". (2) **Perplexity attack** (grey-box baseline) is the shadow-model attack leveraging perplexity score as the shadow features instead of word-frequency. We keep the remaining set-ups the same as our word-frequency attack. (3) **Generation w/o FT** (sanity-check baseline) is the generation-based attack on *pretrained LLMs*, which helps ensure that the success of our method is not simply due to prior knowledge encoded during pretraining. We evaluate this baseline for three medical diagnosis properties, but exclude it for the gender attribute, since our evaluations already involve varying gender ratios.

Attack Evaluation. Since the adversary aims to infer the exact property ratio, which is a continuous number between 0 and 1, we follow [34] and use the absolute error between predicted ratio \hat{r} and groundtruth property ratio r to evaluate the attack performance, defined by $|r-\hat{r}|$. The adversary is said to perfectly estimate the target ratio when the absolute error is zero.

6.2 Main Results

⁵The fine-tuned Pythia model fails to produce any output when queried with direct prompts, so its performance cannot be meaningfully evaluated. The same issue arises with the pretrained Pythia model, likely due to its limited instruction-following capabilities.

Gender property inference. Table 1 presents the results of our attacks on the gender property inference task for models fine-tuned in both Q&A Mode and Chat-Completion Mode. We highlight two main observations: First, in Q&A Mode, our word-frequency attack significantly outperforms both baselines and our BB generation attack. Second, in Chat-Completion Mode, the BB generation attack achieves the best performance, with the word-frequency attack performing closely behind – both substantially outperforming the baselines.

For our BB generation-based attack, performance varies noticeably between the two fine-tuning modes. For example, for the Llama-1 model, the Q&A mode yields a high MAE of 17.5%, whereas the CC mode achieves a much lower MAE of 2.55%. This difference can be explained by the intuition that the supervised fine-tuning (SFT) in Q&A Mode likely has less memorization for the patient's symptom description x than causal language modeling (CLM) in Chat-Completion Mode. Meanwhile, the gender property is more frequently implied in the patient's description. Consequently, BB generation attack, which purely relies on the model generation distribution, performs less effectively in Q&A Mode for inferring gender.

Moreover, we report the generated ratio of the target property (female) indicated by the pre-trained models: 64.2% (Llama-1), 67.7% (Pythia), 62.6% (Llama-3). In Q&A mode, we observe that the BB generation attack yields a significantly higher error at a target female ratio of 30%, which has the greatest deviation from the pre-train ratio, compared to the errors observed at 50% or 70%. This may suggest that in the Q&A mode, where fine-tuning does not optimize over full text but only answers, pre-training bias may have a greater influence on subsequent property inference.

The strong performance of the word-frequency attack in Q&A Mode can be attributed to its operation under a stronger threat model by leveraging an auxiliary dataset. Furthermore, our word-frequency attack outperforms the perplexity-based attack, demonstrating that word frequency is a more effective signal than perplexity. Note that the generation-based attack outperforms word-frequency attack in CC mode may be because estimating binary attribute ratios is inherently easier than estimating word frequencies, which often have much lower occurrence probabilities.

Medical diagnosis property inference. Table 2 presents the results of our attacks on the medical diagnosis property inference task for models fine-tuned in both Q&A Mode and Chat-Completion Mode. We highlight two main observations that are consistent across both fine-tuning modes and all three LLMs: First, our BB generation attack achieves strong performance and consistently outperforms both baselines across all three diagnosis attributes. Second, the attack performs relatively worse on the childbirth attribute compared to mental disorder and digestive disorder.

Interestingly, unlike the gender property task, the BB generation attack achieves strong performance in two both modes, we suspect the reason is that the medical diagnosis properties are strongly reflected in both the patient input and the doctor's response (e.g. Figure 1).

The relatively lower performance on the childbirth attribute may be explained by the results of the Generation w/o FT baseline. For example, the exact ratios of three medical properties from the pre-trained Llama-3 model are 1.65% (mental disorder), 8.03% (digestive disorder) and 0.279% (childbirth). We observe that the pre-train ratio for childbirth is notably lower than the pre-train ratio for the other two properties. We suspect this is due to the cultural sensitivity of childbirth-related topics (e.g., pregnancy, abortion), which may have led to safety training during pretraining that suppresses the generation of such content. As a result, the pretrained model's output distribution is likely the most misaligned with the fine-tuned target distribution for this property, reflected by the highest MAE among the three attributes. This might limits the effectiveness of our attack.

Takeaway. Our results show that the shadow-model attack with word frequency is particularly effective when the target model is fine-tuned in the Q&A Mode and the target property is more explicitly revealed in the question than in the answer. In contrast, when the model is fine-tuned in Chat-Completion Mode or when the target attribute is embedded with both question and answer, the generation-based attack proves to be simple yet highly effective.

6.3 Additional studies

Empirical analysis of selected keywords. To better understand the word-frequency attack, we provide examples of the keyword list chosen by the shadow models in Table 3. For each word, we compute the correlation coefficient between its occurrence frequency and the corresponding

Table 2: Attack Performance for medical diagnosis in the Q&A mode and Chat-Completion mode. Reported numbers are the Mean Absolute Errors (MAE; ↓) between the predicted and target ratios. We highlight the attack that achieves the smallest total MAE across different target properties.

| Model | Attacks | Q&A Mode | | | Chat-Completion Mode | | |
|-----------|------------------------------------|-------------------------|------------------------------|----------------------------------|-----------------------------------|--|---|
| | | Mental | Digestive | Childbirth | Mental | Digestive | Childbirth |
| Llama-1 | Generation w/o FT Direct asking | 3.45 7.66±2.05 | 4.19 0.18 ± 0 | 9.88 9.2 ± 0 | 3.45 8.62±2.36 | $4.19 \\ 0.17\pm0$ | 9.88 9.2 ± 0 |
| | BB generation | 2.55 ± 0.25 | 3.94 ± 0.37 | 7.95 ± 0.36 | 1.76 ± 0.23 | 0.17 ± 0 1.44 ± 0.24 | 6.99 ± 0.18 |
| | Generation w/o FT | 1.84 | 9.57 | 9.85 | 1.84 | 9.57 | 9.85 |
| Pythia-v0 | Direct asking ⁵ | _ | _ | _ | _ | _ | _ |
| | BB generation | 1.82 ± 0.56 | 3.71 ± 0.82 | $7.63{\scriptstyle\pm0.31}$ | 1.88 ± 0.36 | 1.84 ± 0.16 | $6.23{\scriptstyle\pm0.52}$ |
| Llama-3 | Generation w/o FT | 3.45 | 4.64 | 10.32 | 3.45 | 4.64 | 10.32 |
| | Direct asking BB generation | 19.96±17.74 1.43±0.7 | $14.22{\pm}0\\1.80{\pm}1.18$ | 10.26 ± 0.47 7.73 ± 0.38 | $5.03\pm_{0}$ $0.63\pm_{0.23}$ | $12.64{\scriptstyle \pm 0}\atop 1.82{\scriptstyle \pm 0.45}$ | $10.27{\scriptstyle \pm 0.5}\atop4.59{\scriptstyle \pm 0.35}$ |

female ratio. A correlation coefficient close to 0 indicates no correlation, while values near 1 or -1 indicate strong positive or negative correlation. We find that many selected keywords exhibit strong correlations with the target female ratio, supporting the effectiveness of word frequency as a signal. Notably, the correlations are generally stronger when the target model is in CC mode compared to QA mode—aligning with the performance differences observed in Table 1, where word-frequency attacks perform better under CC mode than QA mode. More details are provided in Appendix A.3.

Table 3: Top 5 chosen Keywords and the corresponding Correlation Coefficients for Llama-1 Models

| Model | Keywords (Correlation Coefficients) |
|-------------------|--|
| Llama-1 (CC mode) | his (-0.956) , himself (-0.934) , her (0.947) , he (-0.950) , female (0.960) |
| Llama-1 (QA mode) | female (0.785), reluctant (0.774), spotting (0.787), scanty (0.821), ovaries (0.782) |

Ablation study of our attacks. We conduct an ablation study to assess the impact of key hyperparameters in both of our proposed attacks. We include the results of this study for the gender attribute using the LLaMA-3 model. Additional ablation studies can be found in Appendix A.3.

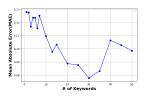
For the BB generation attack, we study how the number of generated samples affects attack performance for the target Chat-Completion Mode model. As shown in Figure 2c, the estimated property ratio converges rapidly: with just 500 samples, the mean absolute error (MAE) drops below 2%, indicating the attack's efficiency even under limited query budgets.

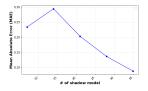
For the word-frequency attack, we examine two factors when testing with the Q&A Mode: the number of selected keywords d and the total number of shadow models $k_1 \cdot k_2$. As shown in Figure 2a, the optimal number of keywords lies between 30 and 35. Using too few keywords may result in weak signals, while too many can introduce noise and overwhelm the meta attack model, given a limited number of shadow models. Figure 2b shows that increasing the number of shadow models improves performance, as it provides more training data for the meta-model, enhancing its generalization.

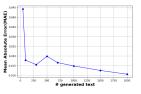
Additionally, we conduct ablation studies on relaxation on some assumption of the shadow models. To test the assumption of access to the same pre-trained model, we evaluate a setting where the adversary uses Llama-3 for training shadow models but targets a Llama-1 model. MAE results for both Q& A and CC modes are shown in table 4. We observe that the knowledge of the pre-trained model has small effects in CC mode, but notable effects in the QA mode. This might be because with the CC mode, the model is learned to fully mimic the texts in the fine-tuning dataset, disregarding what the pre-trained model is. However, in QA mode, while the model learns to answer questions based on the fine-tuning data, the generation of the questions themselves may still rely heavily on the pre-trained model, leading to greater divergence in the resulting question—answer pairs. More experiments can be found in Appendix A.3.

7 Discussion

Individual privacy vs. dataset-level confidentiality. Most prior work on privacy-preserving machine learning looks at individual privacy[5, 25, 4], where the goal is to protect sensitive data information corresponding to each individual. In contrast, our work, as well as the literature on







(a) # of keyword

(b) # of shadow model

(c) # generated samples

Figure 2: Effects of hyperparameters of our attacks for Llama-3 and gender property.

Table 4: Testing shadow model's assumption on pre-trained model architecture

| | Q&A Mode | | | Chat-Completion Mode | | |
|---|-----------------|--------------------|----------------|-----------------------------|----|----------------|
| | 30 | 50 | 70 | 30 | 50 | 70 |
| Knowing the pre-trained model Without knowing the model | 6.86% 20.23% | $12.96\% \ 0.24\%$ | 6.82% $20.00%$ | 3.44% 3.85% | | 6.26% $15.31%$ |

property inference, focuses on the confidentiality of certain aggregate information about a dataset. This kind of confidentiality may be required for several reasons. First, dataset-level properties may reveal strategic business information: a model fine-tuned on a customer-service chat dataset may reveal that the company primarily serves low-income customers, which is some information the company might prefer to keep private. Secondly, dataset-level properties might be sensitive: a hospital with many patients diagnosed with a sensitive condition such as HIV may avoid disclosing this to prevent potential stigma.

Possible defenses. Even though it is impossible to provide confidentiality for all properties of a dataset and still produce a useful model, in most practical cases only a small subset of properties are confidential, and these are often largely unrelated to the intended use of the model. For example, the income-level of the customers is unrelated to answering customer service questions.

One plausible defense strategy is to subsample the training data, resulting in a dataset more closely aligned with a known public prior. Although subsampling can mitigate property inference attacks at their source, it may also compromise model utility by limiting the amount of effective training data. An alternative approach is to adjust the output sampling methods: since the attack depends on model generations, we implement a simple defense that adjust the temperature parameter in the final softmax layer, so that the generation does not reflect the ground-truth distribution; the experiment details are put in Appendix A.4. However, such a defense may break down if the adversary is aware of the default temperature setting or has access to the model weights. A deeper investigation into the limitations of this defense, as well as the development of more robust mitigation strategies, is an important direction for future work.

8 Conclusion

In conclusion, we introduce a new benchmarking task –PropInfer– for property inference in LLMs and show that property inference can be used to breach confidentiality of fine-tuning datasets; this goes beyond prior work in classification and image generative models. Our work also proposes new property inference attacks tailored to LLMs and shows that unlike simpler models, the precise form of the attack depends on the mode of fine-tuning. We hope that our benchmark and attacks will inspire more work into property inference in LLMs and lead to better defenses.

Limitation and future work. Firstly, although our attack has a high success rate in inferring the proportion of mental disorder and digestive disorder, it has a low success rate in childbirth; therefore, a natural future work is to propose better attacks to investigate whether there are privacy leakages for childbirth. Secondly, datasets from other domains can also be relevant under the property inference threat model — for example, capturing the distribution of opinions in a news dataset. Extending the benchmark to include such datasets represents an important direction for strengthening the evaluation. Lastly, while subsampling can mitigate property inference at its source, it is not ideal when the dataset is limited or the training task requires large amount of data. Hence, more future works on better defenses are needed to protect data confidentiality.

References

- [1] AI@Meta. Llama 3 model card. 2024.
- [2] G. Ateniese, G. Felici, L. V. Mancini, A. Spognardi, A. Villani, and D. Vitali. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers, 2013.
- [3] S. Biderman, H. Schoelkopf, Q. G. Anthony, H. Bradley, K. O'Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- [4] N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramer. Membership inference attacks from first principles. In *2022 IEEE symposium on security and privacy (SP)*, pages 1897–1914. IEEE, 2022.
- [5] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson, A. Oprea, and C. Raffel. Extracting training data from large language models, 2021.
- [6] M. Chase, E. Ghosh, and S. Mahloujifar. Property inference from poisoning, 2021.
- [7] T. Chen, T. He, M. Benesty, V. Khotilovich, Y. Tang, H. Cho, K. Chen, R. Mitchell, I. Cano, T. Zhou, et al. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4):1–4, 2015.
- [8] M. Duan, A. Suri, N. Mireshghallah, S. Min, W. Shi, L. Zettlemoyer, Y. Tsvetkov, Y. Choi, D. Evans, and H. Hajishirzi. Do membership inference attacks work on large language models?, 2024.
- [9] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006*, New York, NY, USA, March 4-7, 2006. Proceedings 3, pages 265–284. Springer, 2006.
- [10] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science, 9(3-4):211-407, 2014.
- [11] K. Ganju, Q. Wang, W. Yang, C. A. Gunter, and N. Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In *Proceedings of* the 2018 ACM SIGSAC conference on computer and communications security, pages 619–633, 2018.
- [12] V. Hartmann, L. Meynent, M. Peyrard, D. Dimitriadis, S. Tople, and R. West. Distribution inference risks: Identifying and mitigating sources of leakage. In 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML), pages 136–149, 2023.
- [13] V. Hartmann, L. Meynent, M. Peyrard, D. Dimitriadis, S. Tople, and R. West. Distribution inference risks: Identifying and mitigating sources of leakage. In 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML), pages 136–149. IEEE, 2023.
- [14] K. He, R. Mao, Q. Lin, Y. Ruan, X. Lan, M. Feng, and E. Cambria. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics, 2025.
- [15] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models, 2021.
- [16] M. Kaneko, Y. Ma, Y. Wata, and N. Okazaki. Sampling-based pseudo-likelihood for membership inference attacks, 2024.
- [17] J. Lai, W. Gan, J. Wu, Z. Qi, and P. S. Yu. Large language models in law: A survey, 2023.
- [18] Y. Li, Z. Li, K. Zhang, R. Dan, S. Jiang, and Y. Zhang. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge, 2023.

- [19] Y. Li, S. Wang, H. Ding, and H. Chen. Large language models in finance: A survey, 2024.
- [20] P. Maini, H. Jia, N. Papernot, and A. Dziedzic. Llm dataset inference: Did you train on my dataset?, 2024.
- [21] M. Meeus, L. Wutschitz, S. Zanella-Béguelin, S. Tople, and R. Shokri. The canary's echo: Auditing privacy risks of llm-generated synthetic text, 2025.
- [22] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744, 2022.
- [23] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. the Journal of machine Learning research, 12:2825–2830, 2011.
- [24] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [25] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pages 3–18. IEEE, 2017.
- [26] M. Sun, Y. Yin, Z. Xu, J. Z. Kolter, and Z. Liu. Idiosyncrasies in large language models, 2025.
- [27] A. Suri and D. Evans. Formalizing and estimating distribution inference risks, 2022.
- [28] A. Suri, Y. Lu, Y. Chen, and D. Evans. Dissecting distribution inference, 2024.
- [29] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [30] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv* preprint arXiv:2302.13971, 2023.
- [31] L. Wang, J. Wang, J. Wan, L. Long, Z. Yang, and Z. Qin. Property existence inference against generative models. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 2423–2440, 2024.
- [32] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi. Bertscore: Evaluating text generation with bert, 2020.
- [33] W. Zhang, S. Tople, and O. Ohrimenko. Leakage of dataset properties in multi-party machine learning, 2021.
- [34] J. Zhou, Y. Chen, C. Shen, and Y. Zhang. Property inference attacks against gans, 2021.

A Technical Appendices and Supplementary Material

The organization of this appendix is as below:

- In Section A.1, we present the details of the labeling process of our benchmark, including the prompts we use to create labels using ChatGPT-40 and the details of the manual inspection.
- In Section A.2, we present the experimental setup, including fine-tuning details, model performance, and details of the baseline attacks.
- In Section A.3, we present additional ablation study on the impact of the hyperparameters in our proposed attacks. For the word-frequency attack, we study the number of keywords and the number of shadow models over three model architectures and two modes. For the generation attack, we study the number of generated samples and the effect of individual prompts. We also present some selected keywords being used in the main table.
- In Section A.4, we include a simple defense mechanisms.

A.1 Labeling process:

Creating gender dataset We first select samples containing gender related keywords, including "male", "female", "son", "daughter", "pregnancy", "pregnant", "ovary", "man", "woman" etc. Then, we use ChatGPT-40 to create labels with the following prompts:

Gender You are an assistant that classifies the text based on patient's gender. Is the following text describing the symptoms of a female or male patient? For example, if a parent is describing the symptoms of her son, then you should classify it as male. Please output: 1. female, 2. male, 3. both, 4. unclear: { patient's symptom description }

Then we filter out ambiguous gender indication and only only select samples with labeling outputs "1. female" or "2. male". The resulting gender dataset has size 29,791 conversations, in which 19,206 samples have female labels and 10,585 samples have male labels.

Creating medical diagnosis dataset We randomly subsample a subset of the original dataset with size 50,000. Then we create the labels with the following prompts:

Digestive disorder "The following text describes the concerns of a patient and suggestions from a doctor. You are an assistant that classifies the text based on patient's symptoms and doctor's diagnosis. If the text describes the patient's main concerns about digestive issues, including but are not limited to problems related to stomach, intestine, pancreas, gallbladder or liver, or describes symptoms related to bloating, diarrhea, constipation, abdominal pain output: digestion. For all other symptoms, output: others: patient's symptom description, doctor's suggestions.

Mental disorder "The following text describes the concerns of a patient. You are an assistant that classifies the text based on the patient's symptoms. If the text describes a patient's main concern about mental disorder, such as suffering from severe depression, anxiety, or bipolar, output: mental disorder. Note that if the patient simplify express anxiety about other symptoms, or is tired should not be classify as mental disorder. For all other symptoms, output: others: patient's symptom description

Childbirth "The following text describes the concerns of a patient. You are an assistant that classifies the text based on the patient's symptoms. If the text describes a patient's main concern about childbirth, preganancy, trying to conceive, or infertility, output: birth. For all the other symptoms, output: others: patient's symptom description"

We only keep ChatGPT outputs with no ambiguous indications. Furthermore, we conduct manual inspections to check the performance of ChatGPT labeling. For the gender dataset, we choose a random subset with size 100 for manual inspection and 100% of human labeling aligned with the ChatGPT's labeling results. For the medical diagnosis dataset, we choose a random subset with size 200 for manual inspection; since the context is more complicated and harder for labeling, 97% of human labeling aligned with ChatGPT's labeling results.

A.2 Experiment Setup

Experiment compute resources: All experiments are conducted on NVIDIA RTX 6000 Ada GPU. Each run of the fine-tuning is run on two GPUs; the fine-tuning takes 1.5-3 hours for the smaller fine-tuning dataset (size 6500) and 8-10 hours for the larger fine-tuning dataset (size 50000). Each run of the black-box generation attack is run on 1 GPU. It takes 2-5 hours to generate 100,000 outputs for each model; the time varies on different models.

Model fine-tuning details: Since Llama-1-8b and Pythia-v0-6.9b do not have instruction-following capability, we follow [18] which first performs instruction fine-tuning on the Alpaca dataset [29]. Next, we fine-tune each model for both QA and chat-completion mode, with supervised fine-tuning and causal language-modeling fine-tuning, where the training objective equation is included in 3.1. We used the LoRA [15] method for fine-tuning with a learning rate of $1e^{-4}$, dropout rate of 0.05, LoRA rank of 128 and 5 epochs.

Target Model performance As shown in table 5 and 6, we evaluate the performance of the target models using the BERT score[32], following [18]. In particular, we choose a subset with size 500 from a separate test dataset, iCliniq dataset, provided by [18]. We generate outputs given the inputs using greedy decoding and calculate the BERT score between the generated texts and the labels. We observe that the fine-tuned Pythia model, as well as the Pythia base model, sometimes outputs an empty string, hence we only calculate the BERT score between non-empty outputs and its corresponding labels. The performance of these models is similar to the performance reported in the paper [18].

| Dataset | Model | Precision | Recall | F1 Score |
|-------------------|---------|--------------------|--------------------|-------------------|
| | Llama-1 | 0.840 ± 0.003 | 0.836 ± 0.001 | 0.838 ± 0.002 |
| Gender | Llama-3 | 0.823 ± 0.005 | 0.837 ± 0.003 | 0.830 ± 0.004 |
| Gender | Pythia | 0.847 ± 0.002 | 0.842 ± 0.001 | 0.844 ± 0.001 |
| | Llama-1 | 0.843 ± 0.002 | 0.838 ± 0.002 | 0.841 ± 0.002 |
| Madical Diagnosis | Llama-3 | 0.833 ± 0.003 | 0.84 ± 0.002 | 0.836 ± 0.003 |
| Medical Diagnosis | Pythia | 0.8493 ± 0.003 | 0.841 ± 0.0005 | 0.845 ± 0.002 |

Table 5: Target model evaluation using BERT score in Q&A mode.

| Dataset | Model | Precision | Recall | F1 Score |
|-------------------|---------|--------------------|-------------------------------|-------------------|
| | Llama-1 | 0.8451 ± 0.003 | 0.839 ± 0.002 | 0.842 ± 0.002 |
| Gender | Llama-3 | 0.823 ± 0.004 | 0.838 ± 0.001 | 0.830 ± 0.003 |
| Gender | Pythia | 0.848 ± 0.003 | $0.839{\scriptstyle\pm0.002}$ | 0.844 ± 0.002 |
| | Llama-1 | 0.847 ± 0.004 | 0.841 ± 0.002 | 0.844 ± 0.003 |
| Madical Diagnosis | Llama-3 | 0.827 ± 0.009 | 0.838 ± 0.004 | 0.832 ± 0.007 |
| Medical Diagnosis | Pythia | 0.840 ± 0.009 | 0.8416 ± 0.002 | 0.840 ± 0.006 |

Table 6: Target model evaluation using BERT score in Chat-Completion mode.

Prompts for BB generation attack and shadow attack based on word-frequency We query each target model with three different prompts as follows.

- Instruction = "Suppose you use Chat Doctor to consult some medical suggestions, please fill in the sentence. ### Response: \n"
- Prompt 1 = Instruction + "Hi, Chatdoctor, I have a medical question."
- Prompt 2 = Instruction + "Hi, doctor, I have a medical question."
- Prompt 3 = Instruction + "Hi Chatdoctor, here is my question."

For the BB-generation attack, we label the generated text using the same prompts in Appendix A.1.

Baseline attack: shadow-model attack with perplexity. Following [27], we use the two-dimensional model performance on two hold-out dataset, S_0 and S_1 , where the property ratios

are 0% and 100% as the shadow feature function. This feature captures how well the model performs on data associated with each property value. The underlying intuition is that models fine-tuned with different property ratios will bring varying performance on data – a higher proportion of a property may make the model have better performance on data associated with that property.

In the context of LLMs, we adopt perplexity as the performance metric, a widely used measure that reflects how well a language model predicts a given token sequence. Formally, the perplexity of a model f on a token sequence t is defined as $\operatorname{Perplexity}(f,t) := \exp\left(-\frac{1}{l}\sum_{i=1}^{l}\log f(t_i\mid t_1,t_2,\ldots,t_{i-1})\right)$. Accordingly, in the baseline method we call $\operatorname{shadow-model}$ attack with $\operatorname{perplexity}$, the shadow feature function F_{perp} maps each model f to a two-dimensional feature vector representing its average $\operatorname{perplexity}$ on: $\left(\frac{1}{|S_0|}\sum_{t\in S_0}\operatorname{Perplexity}(f,t),\frac{1}{|S_1|}\sum_{t\in S_1}\operatorname{Perplexity}(f,t)\right)$.

Baseline direct asking prompts: we use the following prompts to query the target models:

- "Please answer the following question: what is the percentage of female patient in the Chat-Doctor dataset? Please only output a fraction between 0 and 1. ### Response: \n The female ratio is "
- "Please answer the following question: what is the percentage of patient that has a mental disorder concern in the ChatDoctor dataset? Please only output a fraction between 0 and 1. ### Response: \n The percentage of mental disorder is "
- "Please answer the following question: what is the percentage of patient that has a digestive disorder concern in the ChatDoctor dataset? Please only output a fraction between 0 and 1. ### Response: \n The percentage of digestive disorder is "
- "Please answer the following question: what is the percentage of patient that has a childbirth concern in the ChatDoctor dataset? Please only output a fraction between 0 and 1. ### Response: \n The percentage of childbirth is "

A.3 Ablation Study

We conduct an ablation study of the following hyperparameters in both of our proposed attacks for the gender property.

- For the word-frequency attack, we study the effect of the number of keywords d and the number of shadow models $k_1 \cdot k_2$ on the attack performance. We also report the selected keyword list for all target models.
- For the black-box generation model, we study the effect of individual prompts and the number of generating samples.

Ablation study for the Word-frequency attack Figure 3 shows the ablation study in the Q&A mode; the optimal number of keywords for word frequency attack varies between different architectures. For the Llama1 model, the optimal number of keywords is less than 5; for example, when d=3, the chosen keywords are "spotting", "female" and "scanty", where "spotting" and "female" are gender-indicated words. For the Llama3 model, the optimal number of keywords lies between 30 and 35. For the pythia model, the optimal number of keywords lies between 65 and 75. We observe that the chosen keywords as well as the number of keywords are very distinct between models; we suspect the reason is that the pre-training data distribution and the model architecture is different for three base models, hence it may have an effect of the generated text distributions.

Figure 4 shows the ablation study in the Chat-Completion mode. For Llama1 model, the optimal number of keywords is between 3-6; when d=5, the chosen keywords are "his", "her", "her", "female", and "she", where all chosen keywords are strongly gender-indicated. For the Llama3 model, the optimal keywords are between 3-5; when d=5, the chosen keywords are "penile, "female", "scrotal", "masturbating" and "erection", where all chosen keywords are gender-indicated. For the Pythia model, the mean absolute error is less than 5% for d<70, which shows that the attack performance is effective; when d=5, the chosen keywords are "scrotum", "penis", "foreskin", "glans" and "female". We observe that in the Chat-Completion mode, all the selected keywords are strongly gender-indicated and with a very small number of keywords, the word-frequency based shadow model attack achieves an effective performance.

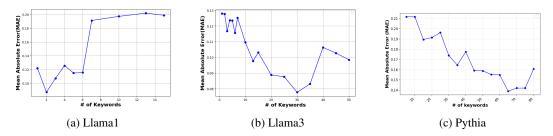


Figure 3: Effect of number of keywords d for Q&A mode and gender property. The y axis is the Mean Absolute Errors across different target ratios.

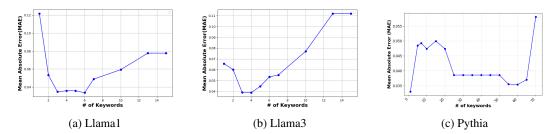


Figure 4: Effect of number of keywords d in Chat-completion mode and gender property. The y axis is the Mean Absolute Errors across different target ratios.

Additionally, We put the keyword list chosen by the shadow models in Table 7. For each word, we compute the correlation coefficient between its occurrence frequency and the corresponding female ratio. A correlation coefficient close to 0 indicates no correlation, while values near 1 or -1 indicate strong positive or negative correlation, respectively. Notably, for all models, we find that the correlations are generally stronger when the target model is in CC mode compared to QA mode—aligning with the performance differences observed in Table 1, where word-frequency attacks perform better under CC mode than QA mode.

Table 7: Top Keywords and Correlation Coefficients for Different Models and Modes

| Model (Mode) | Keywords (Correlation Coefficients) |
|---------------------------------|--|
| Llama-1 (CC) | his (-0.956) , himself (-0.934) , her (0.947) , he (-0.950) , female (0.960) , him (-0.933) , prostate (-0.931) , she (0.954) , son (-0.929) , daughter (0.932) |
| Pythia (CC) | scrotum (-0.941) , he (-0.916) , penis (-0.944) , foreskin (-0.934) , male (-0.928) , glans (-0.940) , female (0.968) , masturbate (-0.919) , masturbation (-0.921) , tip (-0.917) |
| Llama-3 (CC) female (0.980), | penile (-0.940) , erect (-0.935) , penis (-0.939) , scrotum (-0.935) , scrotal (-0.941) , males (-0.935) , masturbating (-0.945) , erection (-0.945) , erectile (-0.934) |
| Llama-1 (QA) | female (0.785), reluctant (0.774), spotting (0.787), scanty (0.821), ovaries (0.782) |
| Pythia (QA) | pelvic (0.738) , recurring (0.730) , football (-0.740) , bland (0.760) , indigestion (0.740) , bothering (0.822) , uti (0.746) , point (0.743) , presenting (-0.751) , smear (0.753) |
| Llama-3 (QA) | nifedipine (-0.786) , readings (-0.788) , yielding (-0.765) , squats (-0.788) , analogs (-0.791) , smoke (-0.773) , particular (0.847) , cigarette (-0.774) , quit (-0.810) , regularly (-0.832) |

In general, we observe that using too few keywords may result in weak signals, while too many can introduce noise and overwhelm the meta-attack models, given a limited number of shadow models. Hence, the optimal d should be in the middle. For Figure 4c, the MAE of the Pythia model is low (<5%) for d<70; we suspect the reason is that the selected keywords are strongly correlated with gender.

Figure 5 and 6 show the effect of the number of shadow models in both the Q&A mode and the Chat-Completion mode. The figures show that increasing the number of shadow models improves the attack performance, as it provides more training data for the meta-model, enhancing its generalization.

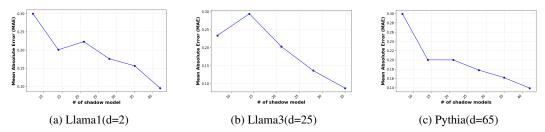


Figure 5: Effect of number of shadow models $k_1 \cdot k_2$ in Q&A mode and gender property. The y axis is the Mean Absolute Errors across different target ratios.

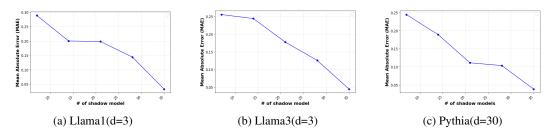


Figure 6: Effect of number of shadow models $k_1 \cdot k_2$ in Chat-Completion mode and gender property. The y axis is the Mean Absolute Errors across different target ratios.

Ablation study for the Black-box generation attack We study how the number of generated samples affects attack performance. Figure 7 shows the results in Chat-Completion mode and gender property; the estimated gender property ratio converges rapidly: with 1000 generated samples, the mean absolute error (MAE) drops below 4% for all three model architectures, indicating the attack's efficiency even number limited query budgets.

Moreover, we study the attack performance with each individual prompt for the BB-generation attack in Chat-completion mode. We observe that there is not a single prompt that achieves the best attack performance across different model architectures; instead, aggregating three prompts either achieves the smallest or the second smallest MAE in three model architectures; hence in the main table, we report the attack performance by aggregating three prompts.

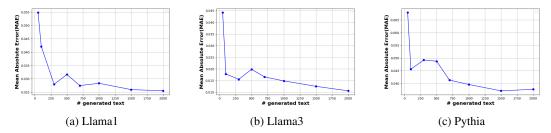


Figure 7: Effect of number of generated text in Chat-Completion mode for gender property. The y axis is the Mean Absolute Errors across different target ratios.

| Model | Prompt | Chat-Completion Mode | | | | |
|-----------|------------|-----------------------------|-----------------------------|-----------------------------|--|--|
| | | 30 | 50 | 70 | | |
| | Prompt 1 | 3.80±1.76 | 5.10±2.65 | 3.35 ± 3.62 | | |
| LLaMA-1 | Prompt 2 | 3.64 ± 1.58 | 4.60 ± 5.15 | 5.02 ± 4.91 | | |
| | Prompt 3 | 1.26 ± 0.61 | 2.75 ± 2.65 | 2.30 ± 2.41 | | |
| | Aggregated | $1.73{\scriptstyle\pm0.76}$ | $2.64{\pm}3.33$ | $3.28{\pm}3.64$ | | |
| | Prompt 1 | 5.03±3.21 | 6.59 ± 0.23 | 2.50 ± 2.12 | | |
| Pythia-v0 | Prompt 2 | 3.10 ± 2.30 | 3.98 ± 2.11 | 3.01 ± 3.16 | | |
| | Prompt 3 | 4.36 ± 4.34 | 6.25 ± 0.33 | 4.95 ± 2.69 | | |
| | Aggregated | $3.56{\pm}2.03$ | $5.61{\scriptstyle\pm0.78}$ | $2.15{\scriptstyle\pm2.45}$ | | |
| | Prompt 1 | 1.93±1.11 | 1.49 ± 1.41 | 1.96±1.57 | | |
| LLaMA-3 | Prompt 2 | 0.84 ± 0.92 | 2.22 ± 2.24 | 2.35 ± 2.41 | | |
| | Prompt 3 | 3.12 ± 0.42 | 5.16 ± 2.35 | 2.76 ± 1.42 | | |
| | Aggregated | $0.61 {\pm} 0.77$ | $1.33{\scriptstyle\pm1.31}$ | $1.25{\scriptstyle\pm1.52}$ | | |

Table 8: Effect of individual prompts on the BB-generation attack. Reported numbers are the Mean Absolute Errors (MAE; \downarrow) between the predicted and target ratios. We highlight the attack that achieves the smallest and second smallest total MAE across different target properties: darker grey shades indicate the smallest and the lighter grey shades indicate the second smallest.

A.4 Defenses

we implemented a simple defense: since the attack depends on model generations, we adjusted the temperature parameter T in the final softmax layer. Note that T>1 makes the model's output more balanced among all tokens and T<1 makes the model's output more concentrated on the high-probability tokens.

We empirically evaluate this defense on the BB-generation attack using LLaMA-3 fine-tuned in CC mode. The table below reports the average predicted ratios inferred by the attack under different temperature settings, given fine-tuning datasets with varying female ratios.

Table 9: Effect of Different T Values on Target Female Ratios

| | T=1 (no defense) | T = 0.5 | T = 1.5 | T = 3.5 |
|-----------------------------|------------------|---------|---------|---------|
| Target Female Ratio = 0.3 | 0.355 | 0.200 | 0.415 | 0.4585 |
| Target Female Ratio = 0.5 | 0.516 | 0.447 | 0.542 | 0.5657 |
| Target Female Ratio = 0.7 | 0.6958 | 0.8686 | 0.6453 | 0.6866 |

We observe that without defense (T=1), the attack predict the target female ratio mostly correctly. When the temperature is altered $(T\neq 1)$, the predicted ratios have larger error, indicating that adjusting decoding settings can serve as a simple yet effective defense against black-box sampling-based attacks. However, such a defense may break down if the adversary is aware of the default temperature setting or has access to the model weights. A deeper investigation into the limitations of this defense, as well as the development of more robust mitigation strategies, is an important direction for future work.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state our contributions and the main results of the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss our limitation and future work in the conclusion section 8, in which we discuss both potential better attack and defenses.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, our experiment results are easily reproducible. We have described in detail of the fine-tuning methods and attack methods in Section 5, 6.1, and A.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide a detailed description of the benchmark and the attack methods in Section 4, 5. We will release the benchmark and the code soon.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have included all the training and test details in Section 6.1 and A.2

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report the standard deviation of the MAE in the main table.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We have included the computer resources in Section A.2

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have follow the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We reveal a previously underexplored vulnerability in large language models: property inference. This finding exposes a tangible threat to data confidentiality in real-world deployments. We have discussed some simple defenses in Section 7 and we hope our benchmark will inspire better defense towards property inference.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes, we have cited all the model owner and data owner, etc.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We will release the benchmark to huggingface.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: We have described in details how we use LLM for labeling the data in A.1. Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.