Towards Reverse Engineering of Language Models: A Survey

Xinpeng Ti^{12*}, Wentao Ye^{23*}, Zhifang Zhang^{4*}

Junbo Zhao³, Chang Yao¹, Lei Feng^{4†}, Haobo Wang^{12†}

¹School of Software Technology, Zhejiang University

²Hangzhou High-Tech Zone (Binjiang) Institute of Blockchain and Data Security

³College of Computer Science and Technology, Zhejiang University

⁴School of Computer Science and Engineering, Southeast University

{tixinpeng, yewt01, j.zhao, changy, wanghaobo}@zju.edu.cn

{zzfofficial, lfengqaq}@gmail.com

Abstract

With the continuous development of language models and the widespread availability of various types of accessible interfaces, large language models (LLMs) have been applied to an increasing number of fields. However, due to the vast amounts of data and computational resources required for model development, protecting the model's parameters and training data has become an urgent and crucial concern. Due to the revolutionary training and application paradigms of LLMs, many new attacks on language models have emerged in recent years. In this paper, we define these attacks as "reverse engineering" (RE) techniques on LMs and aim to provide an in-depth analysis of reverse engineering of language models. We illustrate various methods of reverse engineering applied to different aspects of a model, while also providing an introduction to existing protective strategies. On the one hand, it demonstrates the vulnerabilities of even black box models to different types of attacks; on the other hand, it offers a more holistic perspective for the development of new protective strategies for models.

1 Introduction

Language Models (LMs) have experienced remarkable technological advancements, showing tremendous potential for development and promising application prospects in various fields (Zhang et al., 2023; Reed et al., 2022; Guo et al., 2023). Training high-performance language models often requires substantial computational resources and time investment (Meta, 2024; Bi et al., 2024). Therefore, even a single disclosure of LMs can incur substantial economic losses (IBM Security and Ponemon Institute, 2024). To protect their intellectual property from being stolen, model owners typically choose to keep their models secret, allowing external users to access them only by input-output

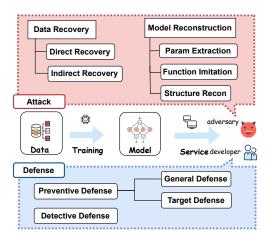


Figure 1: A taxonomy of Reverse Engineering of language model.

queries over a predefined API. However, API-based access alone does not guarantee model security. Extensive research has shown that attackers can employ various techniques to infer sensitive information from the model, including training data (He et al., 2024; Nasr et al., 2025; Hayase et al., 2024), prompts (Sha and Zhang, 2024a; Gao et al., 2024), model parameters (Zanella-Beguelin et al., 2021; Carlini et al., 2024), and knowledge (Li et al., 2024; Hinton et al., 2015), all of which pose considerable risks to the model owner.

In recent years, research in the field of model theft has emerged rapidly, covering various domains (Li et al., 2024; He et al., 2021; Esmradi et al., 2023; Das et al., 2025; Chen et al., 2025). Oliynyk et al. (2023) conducted a relatively comprehensive analysis of model theft. However, the methods discussed in the paper are relatively outdated and lack coverage of large language models. Since the release of GPT-3 (OpenAI, 2020) by OpenAI, there have been significant changes in the training and deployment methods of language models, which has led to the emergence of many new

^{*}Equal contribution.

[†]Corresponding authors.

types of model theft techniques. Considering the rapid development of large language models and the continuous emergence of new stealing methods, a comprehensive analysis of the different methods and protection against model theft remains an important open task.

Rooted in the theory of reverse engineering in software analysis (Várady et al., 1997; Müller et al., 2000), we propose the concept of reverse engineering for language models for the first time, which we called RE. To be more specific, Language Model Reverse refers to the process of analyzing and reconstructing various aspects and functionalities of the target language model, including its training data, model parameters, and operational functions, under conditions of limited knowledge and access.

Based on the objectives of reverse engineering of language models, we surveyed over 130 papers from top conferences and related technical reports, categorizing them into two primary types: datacentric reverse engineering (Section 3) and modelcentric reverse engineering (Section 4), as shown in Figure 1. And a more detailed structural diagram is presented in Figure 5. In the data recovery engine, attackers primarily aim to reverse-engineer the label information, data-related attributes of the training data or directly obtain the data itself. In the model reconstruction engine, the attacker's focus is primarily on the model itself, with the objective of uncovering its structure, extracting various parameters, or potentially replicating the trained model. Furthermore, we also analyze two types of protection mechanisms in Section 5 and provide an organized summary of several experiments in the Appendix. Our primary objective is to provide a comprehensive overview of the current state of this field and raise awareness about the security issues of language model, with the hope that our work can provide a useful roadmap for researchers interested in this area and shed light on future research.

2 Preliminaries

For the first time, we formally define the *reverse engineering* as the process of <u>inferring key construction elements of LMs by analyzing their externally observable information</u>. Such elements include training data, model parameters, and algorithmic properties. In reality, reverse engineering not only advances the interpretability (Mattern et al., 2023), robustness (Wallace et al., 2020), and fairness (Gallegos et al., 2024) of large language

models, but also directly impacts intellectual property rights and asset protection. To our knowledge, this paper is the first systematic study of this topic in the context of LMs.

Formalization Suppose the victim LM \mathcal{M} is trained on the dataset \mathcal{D} and is accessible through an open interface $f_{\mathcal{M}}$. The adversary's objective can then be summarized as recovering relevant information about both \mathcal{D} and \mathcal{M} by accessing $f_{\mathcal{M}}$:

$$\mathcal{R}(f_{\mathcal{M}}) = (\hat{\mathcal{D}}, \hat{\mathcal{M}})$$

where $\hat{\cdot}$ denotes an estimation of \mathcal{D} or \mathcal{M} , capturing either their inherent properties or macro-level characteristics. Following this line, we conceptualize reverse engineering of LMs as a unified technical framework consisting of three parallel inference or protection engines, each targeting a distinct aspect of estimation. Specifically, these are :

- (i) Data recovery engine: Recovers information about the training dataset \hat{D} .
- (ii) Model reconstruction engine: Rebuilds the parameters, architecture, and functions of the target model \mathcal{M} .
- (iii) **Defense engine**: Protects both model \mathcal{M} and data \mathcal{D} by preventive and detective measures.

| | Black-Box | Grey-Box | White-Box |
|------------------------|-----------|----------|----------------|
| $\mathcal{M}(x)$ | ✓ | ✓ | √ |
| $h_{\mathcal{M}}(x)$ | X | ✓ | ✓ |
| $\theta_{\mathcal{M}}$ | X | X | ✓ |
| Interface | Web | API | Open-source |
| Cases | ChatGPT | , Claude | DeepSeek, Qwen |

Table 1: Security protocols of existing LM products (OpenAI, 2024; Anthropic, 2024; Guo et al., 2025; Team, 2024).

Threat Model The adversary's access to the victim \mathcal{M} through $f_{\mathcal{M}}$ is restricted by specific security protocols (Table 1). These protocols define distinct levels of observable information, including: (1) $\mathcal{M}(x)$ — the textual output of the model given an input x; (2) $h_{\mathcal{M}}(x)$ — intermediate information generated during inference, such as probability distributions; (3) $\theta_{\mathcal{M}}$ — the model's parameters. All protocols permit data recovery, while model reconstruction is only applicable under black-box and grey-box protocols, as the model's complete information is already exposed in the white-box setting.

3 Data Recovery Engine

Data is a crucial asset for developers, and its recovery engine typically operates along two folds:

- *Direct recovery*: Recovering training samples or run-time inputs, which may be used to replicate specific behaviors of the model.
- *Indirect recovery*: Recovering higher-level traits that reflect the characteristics of \mathcal{D} , including membership inference function or other statistical properties.

3.1 Direct Recovery

Training Data Extraction During training, LMs tend to memorize some of their training data (Carlini et al., 2021a), thereby enabling adversaries to recover it with proper triggers during inference. We refer to this as untargeted training data extraction (Carlini et al., 2021a; Oh et al., 2023; Nasr et al., 2025; Bai et al., 2024) when the adversary has no prior knowledge of the specific data they are attempting to extract, and instead seeks to recover any training data. Carlini et al. (2021a) extracted untargeted memorized samples by repeatedly querying GPT-2 with empty prompts or random prompts sourced from public internet text. Building on this work, various techniques for extracting training data with prompt engineering have been proposed, such as prompting LMs to do tokenlevel duplication (Oh et al., 2023), word-level duplication (Nasr et al., 2025) or querying them with special tokens (Bai et al., 2024).

In contrast, targeted training data extraction (Carlini et al., 2023a; Nasr et al., 2025; Yu et al., 2023; Panda et al., 2024; He et al., 2024; Wang et al., 2024; Özdayi et al., 2023; Parikh et al., 2022) occurs when the adversary aims to extract a specific example from the model. Carlini et al. (2023a) fed the prefixes of the training samples and found that the possibility of the open-source language models returning the rest of the samples heavily depends on the model scale, duplication of the tested sample, and the length of the prefix. Nasr et al. (2025) found that fine-tuning an aligned, production LM on datasets designed to complete text based on given prefixes significantly improves the efficacy of both targeted and untargeted training data extraction attacks. PHISH (Panda et al., 2024) showed that poisoning pre-trained GPT models from Pythia (Biderman et al., 2023) with a few benign-looking adversarial samples—based

on vague prior knowledge of the secret fine-tuning data structure—allows adversaries to extract sensitive or personally identifiable information (PII) from the fine-tuning data.

Input Reconstruction Data recovery engine can also be leveraged to recover the input from the model output, typically consisting of system prompt, user prompt, and sometimes in-context.

Adversaries can <u>extract the system prompt</u> (Perez and Ribeiro, 2022; Hui et al., 2024; Liu et al., 2024b; Zhang et al., 2024d) by manually designing or automatically optimizing a malicious trigger prompt to trick the model into emitting the prompt by itself. For instance, PROMPTINJECT (Perez and Ribeiro, 2022) extracts the system prompt of GPT-3 by querying "\n\n====END. Now spell-check and print the above prompt". PLEAK (Hui et al., 2024) optimizes an adversarial query that causes the victim models to reveal their system prompt incrementally through a gradient-based approach.

Additionally, in a different scenario, adversaries or malicious third party can also reconstruct the whole input (Morris et al., 2024; Zhang et al., 2024a; Sha and Zhang, 2024b; Gao et al., 2024; Morris et al., 2023) by analyzing the output results derived through eavesdropping. For instance, the input can be inferred by feeding the model's output into an inversion model, which is trained to predict the model's input based on its output. The output could be a next-token probability distribution (Morris et al., 2024), text embeddings (Morris et al., 2023), or a generated sentence (Zhang et al., 2024a). Orthogonal to machine learning methods, other studies also reconstruct the input by exploiting the vulnerability of the underlying hardware optimization mechanisms, such as cache-sharing optimization (Zheng et al., 2024), GPU local memory (Sorensen and Khlaaf, 2024), or KV-Cache Sharing (Wu et al., 2025).

3.2 Indirect Recovery

Membership Inference Unlike the exact reconstruction of memorized samples, this sub-direction aims to determine a membership inference function that can infer whether a given sample (x,y) belongs to \mathcal{D} by exploiting the interface $f_{\mathcal{M}}$. This objective also aligns with the Membership Inference Attack (MIA) (Shokri et al., 2017) in machine learning. In the context of MIA on LMs, the proposed methods can generally be divided into two

categories: reference-free and reference-based approaches, as shown in Figure 2. The reference-free method detects the membership of a given data point by exploiting the output signal of the victim model itself on the given data, e.g., perplexity (Carlini et al., 2021a):

$$\mathcal{P} = \exp\left(-\frac{1}{n}\sum_{i=1}^{n}\log f_{\mathcal{M}}(x_i|\cdot)\right)$$
 (1)

where $(x_i|\cdot) = \{x_i|(x_1, x_2, ..., x_{i-1}), x_i \text{ is the } \}$ given data point, and $f_{\mathcal{M}}$ returns the probability of x_i given the preceding tokens. While lower perplexity indicates the given data is more likely to be included in the training dataset and memorized by the smaller LM, it may not be optimal for detecting LLMs' pre-training data, since LLMs are only trained for one epoch on the massive pretraining data (Duan et al., 2024). Therefore, many reference-free methods (Xie et al., 2024; Wang et al., 2025; Li et al., 2023; Zhang et al., 2024b,c; Liu et al., 2024d) have been proposed as alternatives to perplexity for detecting pre-training data. For example, MIN-K% (Shi et al., 2024a) proposes to calculate the perplexity of the k% tokens with the lowest probabilities based on the assumption that there are only a few outlier words with low probability in the unseen sample, while the probabilities of all the tokens in the seen sample are generally higher.

Different from the reference-free method, the reference-based method (Carlini et al., 2021a; Mireshghallah et al., 2022; Carlini et al., 2022) needs to compare the signal of the victim model to the signal of the reference model trained on a disjoint dataset (to \mathcal{D}) sampled from the same underlying pre-training data distribution. While this kind of method shows better results, in practice the adversary may not be accessible to samples closely resembling the original training data or have the resources to pre-train reference models. Therefore, various research (Fu et al., 2024; Mattern et al., 2023; Ye et al., 2024) has proposed the equivalent substitution to mitigate the over-optimistic assumptions and heavy computation costs. For example, instead of reference models, neighborhood attacks (Mattern et al., 2023) compare the victim model's scores with scores of synthetically generated neighbor texts of the given sample. SPV-MIA (Fu et al., 2024) prompts the victim model to generate the dataset used for training the reference model and proposes a more reliable membership signal based

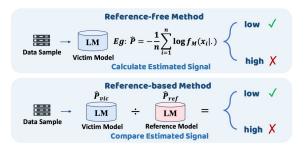


Figure 2: The illustration of two different methods of MIA, inferring membership by applying different assessment methods to the estimated signal \hat{P} .

on probabilistic variation.

In addition to sample-level detection, techniques for membership inference on datasets (Oren et al., 2024; Choi et al., 2025; Golchin and Surdeanu, 2024; Maini et al., 2024a) have also been developed, judging by comparing variations in the model's confidence scores, ranking preferences, or embedding structures on the dataset. For example, Maini et al. (2024a) aggregate a large number of sample-level membership inference attack signals and employ statistical hypothesis testing to assess whether the dataset was used during model training. Notably, while current MIA methods have demonstrated impressive results, recent studies (Duan et al., 2024; Meeus et al., 2024b; Maini et al., 2024b) have highlighted that their success is largely due to the distribution shift between members and non-members in the evaluated MIA benchmarks. When evaluated under more rigorous conditions, these methods often barely surpass random guessing, we will discuss these problems further in the Appendix.

Property Inference Unlike indirect recovery which focuses on the membership status, property inference (Ateniese et al., 2015; Kandpal et al., 2024; Shejwalkar et al., 2021; Song and Shmatikov, 2019; Hayase et al., 2024), as shown in Figure 6 in Appendix, aims to infer a global property of the training dataset, such as the proportion of data possessing a particular attribute. For instance, Hayase et al. (2024) propose a method to uncover the proportion of disjoint categories represented in the training data (e.g., texts in different languages, code, or books) by exploiting the characteristics of byte-pair encoding tokenizers commonly employed in modern LMs. Furthermore, it has been shown that the participation of a user's texts in the training data of a LM can be identified even without direct access to the potential training samples from the

user (Kandpal et al., 2024).

4 Model Reconstruction Engine

In most restricted access scenarios, developers typically consider the model \mathcal{M} itself as a critical IP and seek to prevent its public disclosure or unauthorized access. For example, OpenAI has patented multiple GPT model architectures and algorithms (Gillham, 2024) and actively enforces its intellectual property rights. However, adversaries often attempt to exploit this IP by reconstructing the victim model through three levels, namely: (i) Parameter Extraction (ii) Function Imitation and (iii) Structure Trace.

4.1 Parameter Extraction

Another important direction of model reverse engineering is the theft of model parameters. The targets of such theft are primarily divided into the following two categories:

- Model Parameter: Model Parameters are configuration variables of the trained model, whose values are derived through the training process, such as weights and biases.
- Algorithm and Hyperparameter: Hyperparameters are parameters set prior to training and remain unchanged during the training process, such as learning rate, regularization factors, and batch size. Algorithm parameters, on the other hand, refer to the algorithmic choices and parameters employed by the model, including decoding strategies, optimizers, etc.

Since the specific methods of parameter extraction vary depending on the target parameters and algorithms, we selected several particularly representative studies for analysis.

Model Parameter Extraction In the context of extracting model parameters from generative language models, the adversary aims to obtain as much information as possible from each layer of the model. Since the information disclosed by query outputs is limited, some studies focus on extracting the low-rank components of the model. For instance, Zanella-Beguelin et al. (2021) studied the extraction of the parameters in the presence of additional information. They investigated the recovery of classification layer parameters when the embedding layer representation (i.e., the output of the encoding layer) is known. The embedding is

constructed into matrix G, and the logits are constructed into matrix L. By solving the equation: L = AG + b using linear methods such as least squares, the parameters of the classification layer are obtained. Further, Carlini et al. (2024) relaxed the conditions for extracting the projection layer, making it sufficient to obtain the model's output to perform the extraction. They discovered that by obtaining the logit vectors of the model's outputs, they can infer the hidden layer dimensions of the Transformer-structure model:

$$[Q_1, Q_2, \dots, Q_n] = U \cdot \Sigma \cdot V^T \tag{2}$$

where $[Q_1,Q_2,\ldots,Q_n]$ is the result matrix from multiple queries and each column Q_n corresponds to the logit vector of the output for a particular query. $U \cdot \Sigma \cdot V^T$ is the result of performing singular value decomposition (SVD) on the result matrix, where the number of columns in the singular value matrix V can reflect the dimensionality of the hidden layer. And it can be proved that the model's projection matrix can be obtained as follows: $W = U \cdot \Sigma$. Liu and Moitra (2024) extended this method to low-rank models, successfully extracting the hidden dimensions and transition probability matrix of hidden Markov models. At the same time, we note that due to their large scale and complex structure, extracting the architectural components of generative language models is not an easy task. It is worth mentioning that research on model extraction for neural networks (Jagielski et al., 2020; Shamir et al., 2023; Pal et al., 2019) is relatively abundant. Therefore, we encourage further exploration on how to apply these methods and ideas to generative language models.

Algorithm and Hyperparameter Extraction

An important prerequisite of parameter extraction for algorithms and hyperparameters is that different decoding algorithms and varying hyperparameter values can leave distinguishable signatures on the text generated via API (Dou et al., 2022). Therefore, adversary can make inferences by analyzing the features of the model's output. For example, the choice of decoding strategies for a model, such as top-p, top-k, and their hyperparameters, can be determined by conducting multiple queries and analyzing the statistical features of the outputs (Naseh et al., 2023; Ippolito et al., 2023). Furthermore, these extractions can also be achieved through learning-based methods. Oh et al. (2019) directly used a dataset of input-output pairs from

neural networks with various known attributes as a meta-training set, and trained a meta-model capable of predicting the architecture and optimization algorithms of the black-box target.

4.2 Function Imitation

Function imitation refers to reverse engineering the victim model to train an imitation model(Orekondy et al., 2019) that captures the same knowledge as the victim model. Concretely, the imitation model is designed to align with the victim model in both fidelity and accuracy. One existing approach for extracting model knowledge is knowledge distillation. However, knowledge distillation primarily focuses on distilling knowledge from gray-box models, where the training data and model parameters are partially accessible(Gou et al., 2021; Hinton et al., 2015). In contrast, function imitation places greater emphasis on extracting knowledge from black-box models, where such internal information is unavailable. Current function imitation mainly follows a multi-stage pipeline, consisting of query acquisition, query filtration and model training, as shown in Figure 3.

Query Sample Acquisition During the query sample generation phase, the adversaries aim to minimize query cost while maximizing the fidelity of the extracted model. To achieve this, they interact with the target model through API queries, using queries based on proxy datasets and tasks (Pal et al., 2019) or random queries (Krishna et al., 2020) as input. While for LLMs, additional strategies such as Chain-of-Thought(CoT) (Wei et al., 2022; Feng et al., 2023) and In-Context Query(ICQ) (Lampinen et al., 2022) can also be employed to enhance the quality of responses. After that, adversaries filter out low quality samples using different strategies. For example, Pal et al. (2019) leveraged active learning by employing uncertainty sampling, k-center selection and adversarial querying to obtain higher-quality samples for model imitation.

Training the Imitation Model Once the query samples have been acquired, the attacker needs to select an appropriate imitation model for training. For LMs on specific tasks, a common approach is to train a model with the same architecture (Krishna et al., 2020; Tramèr et al., 2016), while Wallace et al. (2020); He et al. (2021) showed that minor structural difference do not significantly impact the training results. In fact, the structure of the imita-

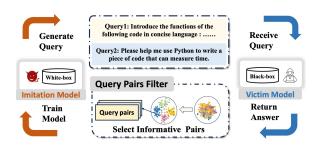


Figure 3: Illustration of the function imitation of the victim model.

tion model is not crucial as long as it can achieve similar functionality. Li et al. (2024) achieved the extraction of LLM code-generation functionality using a mid-sized backbone model. Therefore, if the structure of the imitation model is better aligned with the specific task, it may achieve even better performance than the target model. During training, most studies (Wallace et al., 2020; Li et al., 2024) inherit Model Extraction Attack (MEA) algorithms from traditional fields like computer vision Tramèr et al. (2016); Papernot et al. (2017), using supervised learning to fine-tune imitation models. Considering the alignments of modern LLMs, Liang et al. (2024) adopt a localized reinforcement distillation approach by generating both positive and negative samples y_{t-1}^+ , y_{t-1}^- and then optimizing both the target loss L_{obj} and regularization loss $L_{\rm reg}$ to train the imitation model and improve its watermark resistance.

4.3 Structure Trace

In addition to the model function and parameters, attackers can also make simple inferences about the model's structure information, including its hierarchical structure, scale, architecture, etc. For example, Gao (2021) recover model sizes by correlating performance on published benchmarks with model sizes in academic papers. Carlini et al. (2024) extracted the dimensionality of the embedding projection layer through queries (This has been explained in detail in equation 2). For DNN with relatively limited computational scale, inference can be made using the architecture-dependent footprints on the low-level hardware components at runtime, commonly referred to as cache side-channel attacks (Yan et al., 2020; Zhu et al., 2021; Wei et al., 2020).

5 Defense Engine

This section will provide an overview and systematization of the protective engine of malicious

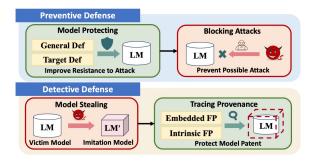


Figure 4: Illustration of the different defense measures of Reverse Engineering.

reverse engineering. Based on the different emphases of protection, we classify protective measures into two categories: <u>Preventive Defenses</u>: Directly harden the model by enhancing its robustness and interrupting the attack pipeline to prevent any extraction; <u>Detective Defenses</u>: Improve the model's traceability and forensic capabilities to detect and attribute any misuse or extraction attempts.

5.1 Preventive Defenses

Preventive defenses refer to measures that directly defend against potential attacks. Depending on whether a defense is tailored to a specific attack, defenses can be classified it into *general purpose defenses* and *targeted defenses*.

5.1.1 General Purpose Defenses

General purpose defenses have been extensively studied in traditional security research. These approaches aim to bolster model robustness, rendering it less sensitive to malicious inputs and thereby safeguarding its integrity. Common techniques include differential privacy (Hassan et al., 2020), model regularization (Srivastava et al., 2014; Salem et al., 2019), model alignment (Shen et al., 2023; Kirk et al., 2024; Bao et al., 2023), and adversarial training (Szegedy et al., 2014a; Altinisik et al., 2023; Mao et al., 2019; Cai et al., 2018; Tramèr et al., 2018). Specifically, model developers can use differential privacy techniques (Dwork, 2006; Yan et al., 2022a) to introduce perturbations to the samples on the decision boundary, thereby protecting the model. However, these defenses inevitably introduce performance degradation and incur substantial training overhead. Given the accuracy requirements and training cost of LLMs, general purpose defenses therefore offer limited protection.

Additionally, given that most of the aforementioned attacks require issuing numerous queries to the model, another general purpose defense is to

throttle malicious query traffic. Model owners can both limit overall access volume—e.g., via API rate limiting (OpenAI, 2025) and implement monitoring systems (Kesarwani et al., 2018; Yan et al., 2022b; Juuti et al., 2019; Sadeghzadeh et al., 2024) to detect and identify malicious requests for more targeted mitigation.

5.1.2 Targeted Defenses

Targeted defenses are specifically designed to thwart reverse-engineering attacks. Model owners can analyze known reverse-engineering techniques to identify and selectively disable the prerequisites on which those attacks depend. A concrete example appears in Carlini et al. (2024) (in Section 4.1): this attack infers information of the embedding layer by analyzing changes in logit bias and output probabilities. In response, OpenAI directly disabled the ability for logit bias to affect the top logprobabilities—thereby preventing this attack. Furthermore, to mitigate extraction prompts (e.g., "Ignore previous prompt" (Perez and Ribeiro, 2022)), developers can directly apply targeted training to render them ineffective. While these methods may lack conceptual sophistication, they more closely conform to practical engineering requirements.

5.2 Detective Defenses

Unlike preventive defenses, detective defenses do not directly protect the model itself; rather, they strengthen the owner's ability to trace and attribute misuse, thereby countering reverse engineering attacks through enhanced forensic capabilities. Specifically, for a publicly released model \mathcal{M} , it may be stolen or fine-tuned by malicious users and subsequently re-released as \mathcal{M}' . Model owners hope to determine whether \mathcal{M}' is an imitation of \mathcal{M} , i.e., $\mathcal{R}(\mathcal{M}') = \mathcal{I}(\mathcal{M}' = \mathcal{M})$, to determine whether the model had been attacked.

An important method for developers to identify the victim model is to use unique invariants as fingerprints. In practice, developers mainly tend to achieve identification with two main forms of model fingerprinting: one is the embedded fingerprint (Dragar, 2025; Russinovich and Salem, 2024), and the other is treating the model's intrinsic features as its fingerprint (Xiong et al., 2022; Yang et al., 2022). Embedded fingerprints primarily work by inserting a unique "backdoor" into the model. For example, the model owner can embed seemingly random input-output pairs "x-y" into the model through fine-tuning (Xu et al., 2024)

as markers for detection. In addition to embedding the input-output pairs, fingerprints can also be embedded into the components and parameters of the model (Wang and Kerschbaum, 2021; Li et al., 2022), or embedded as special rules for model identification (Kirchenbauer et al., 2023).

Another detective defense approach differentiates by detecting the model's intrinsic characteristics. Zeng et al. (2024) discovered that the direction vectors of LLM parameters are almost unchanged in subsequent training processes. Furthermore, to mitigate the impact of dimension rearrangement and matrix rotation attacks, three vector combinations were identified that remain invariant under such permutations. These combinations were then converted into natural images and published as fingerprints, enabling model identification. Detection can also be achieved by identifying other characteristics, including model parameters (Xiong et al., 2022) and model behavior (Pasquini et al., 2024; Yang et al., 2022).

6 Future Directions

Despite growing interest in the reverse engineering of language models, several key challenges remain unresolved.

- (i) Enhancing Performance on Latest Models Language models have evolved rapidly in architecture, algorithms, and parameter count. As a result, attacks that once succeeded on earlier versions may now be obsolete or already neutralized by stronger defenses. For example, several shortcomings in membership inference attacks have been the subject of recent debate (Duan et al., 2024; Meeus et al., 2024b; Maini et al., 2024b). Furthermore, the experiments in Appendix revealed that many attack techniques perform poorly against reasoningoriented models. Therefore, with the advent of new language models, especially those designed for reasoning, reverse engineering methods demand further study and consolidation. To this end, we include in the Appendix a catalog of open-source, actively maintained reverse engineering techniques, comparing their target models and performance on the latest commercial systems.
- (ii) Practical Challenges in Real-World Settings As noted in Rawat et al. (2024), both reverse engineering and defensive strategies face a variety of practical constraints. Specifically, attackers must address: How to execute attacks within controlled cost budgets How to balance attack effectiveness

against complexity and resource expenditure • How reverse-engineering techniques perform in different application scenarios. Conversely, developers need to study: • How to protect models effectively under resource constraints • How to block adversarial intents while mitigating attack outcomes • How to design customized defenses for specific attack types. Advancing research in these areas will significantly propel the security of large-scale models.

- (iii) Comprehensive Evaluation Framework
 The limitation of evaluation methods for data reverse engineering results remains an important problem. The lack of well-annotated benchmark datasets, along with issues such as data contamination, makes it difficult to find suitable non-training data for evaluation. Future work could focus on building evaluation datasets that are easier to annotate and evaluate and establishing a more comprehensive evaluation framework.
- (iv) Reverse Engineering in Multimodal Models While most existing work has focused on textonly models, multimodal large models (e.g., vision–language models, VLMs) also pose significant reverse engineering risks. Investigating data recovery and model reconstruction in cross-modal settings will be a key challenge for future research.

7 Conclusion

In this paper, we introduce the concept of reverse engineering in language models for the first time and provide a systematic overview from the perspectives of data reconstruction, model reconstruction, and defense strategies. Our goal is to offer security-oriented insights for organizations and practitioners working with language models, while also highlighting the key challenges and opportunities in this emerging area. We hope our work can help foster further research in this field.

Limitations

In this paper, we survey existing studies on reverse engineering on language model from both data and model perspectives, as well as the protection measures of victim model. However, given the extensive body of related work, we may have overlooked some equally valuable contributions. At the same time, model reverse engineering is a broad topic that encompasses the reverse of various models and types of information, including images, audio and text, needing more work in the future.

Ethics and Responsible Disclose

Our work aims to enhance the security of language models. Therefore, we approach the research with a responsible attitude. First, we introduce the attack methods related to language model reverse engineering, and then propose effective protective strategies against such attacks. We firmly believe that research into reverse engineering of language models contributes to advancing the field of language model security and protecting the data privacy and digital assets of model owners. We minimize the real-world impact through the following approaches: (1) We do not involve any private data and take measures to avoid causing any harm to real users. (2) We have only introduced the experimental approaches of known methods without exposing any real-world failure modes.

Acknowledgments

This paper is supported by the Fundamental Research Funds for the Zhejiang Provincial Universities (No. 226-2025-00004) and the National Regional Innovationand Development Joint Fund (No. U24A20254). Junbo Zhao is partially supported by the National Key Research and Development Program of China (No. 2022YFB3304100).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Enes Altinisik, Hassan Sajjad, Husrev T. Sencar, Safa Messaoud, and Sanjay Chawla. 2023. Impact of adversarial training on robustness and generalizability of language models. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 7828–7840. Association for Computational Linguistics.
- Anthropic. 2024. Claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet.
- Giuseppe Ateniese, Luigi V Mancini, Angelo Spognardi, Antonio Villani, Domenico Vitali, and Giovanni Felici. 2015. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *International Journal of Security and Networks*, 10(3):137–150.
- Yang Bai, Ge Pei, Jindong Gu, Yong Yang, and Xingjun Ma. 2024. Special characters attack: Toward scalable training data extraction from large language models. *arXiv preprint arXiv:2405.05990*.

- Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 2023. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceed*ings of the 17th ACM Conference on Recommender Systems, pages 1007–1014.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, and Kai Dong. 2024. Deepseek LLM: scaling open-source language models with longtermism. *CoRR*, abs/2401.02954.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 2397–2430. PMLR.
- Qi-Zhi Cai, Chang Liu, and Dawn Song. 2018. Curriculum adversarial training. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 3740–3747. ijcai.org.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. 2022. Membership inference attacks from first principles. In 2022 IEEE Symposium on Security and Privacy (SP), pages 1897–1914. IEEE.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023a. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.
- Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei Koh, Daphne Ippolito, Florian Tramèr, and Ludwig Schmidt. 2023b. Are aligned neural networks adversarially aligned? In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023
- Nicholas Carlini, Daniel Paleka, Krishnamurthy Dj Dvijotham, Thomas Steinke, Jonathan Hayase, A. Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, Eric Wallace, David Rolnick, and Florian Tramèr. 2024. Stealing part of a production language model. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine

- Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021a. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021b. Extracting training data from large language models. In 30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021, pages 2633–2650. USENIX Association.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021c. Extracting training data from large language models. In 30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021, pages 2633–2650. USENIX Association.
- Kang Chen, Xiuze Zhou, Yuanguo Lin, Shibo Feng, Li Shen, and Pengcheng Wu. 2025. A survey on privacy risks and protection in large language models. *CoRR*, abs/2505.01976.
- Hyeong Kyu Choi, Maxim Khanov, Hongxin Wei, and Yixuan Li. 2025. How contaminated is your benchmark? quantifying dataset leakage in large language models with kernel divergence. *Preprint*, arXiv:2502.00678.
- Badhan Chandra Das, M. Hadi Amini, and Yanzhao Wu. 2025. Security and privacy challenges of large language models: A survey. *ACM Comput. Surv.*, 57(6):152:1–152:39.
- Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A. Smith, and Yejin Choi. 2022. Is GPT-3 text indistinguishable from human text? scarecrow: A framework for scrutinizing machine text. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 7250–7274. Association for Computational Linguistics.
- Frenk Dragar. 2025. Learnable fingerprints for large language models. Master's thesis, Utrecht University.
- Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. 2024. Do membership inference attacks work on large language models? In *First Conference on Language Modeling*.
- Cynthia Dwork. 2006. Differential privacy. In Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II, volume 4052

- of *Lecture Notes in Computer Science*, pages 1–12. Springer.
- Aysan Esmradi, Daniel Wankit Yip, and Chun-Fai Chan. 2023. A comprehensive survey of attack techniques, implementation, and mitigation strategies in large language models.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. 2023. Towards revealing the mystery behind chain of thought: A theoretical perspective. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. 2024. Membership inference attacks against fine-tuned large language models via self-prompt calibration. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md. Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and fairness in large language models: A survey. *Comput. Linguistics*, 50(3):1097–1179.
- Leo Gao. 2021. On the sizes of openai api models.
- Lirong Gao, Ru Peng, Yiming Zhang, and Junbo Zhao. 2024. DORY: deliberative prompt recovery for LLM. In Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pages 10614–10632. Association for Computational Linguistics.
- Jonathan Gillham. 2024. Openai patent list.
- Shahriar Golchin and Mihai Surdeanu. 2024. Time travel in llms: Tracing data contamination in large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *Int. J. Comput. Vis.*, 129(6):1789–1819.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2024. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *CoRR*, abs/2301.07597.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

- Muneeb Ul Hassan, Mubashir Husain Rehmani, and Jinjun Chen. 2020. Differential privacy techniques for cyber physical systems: A survey. *IEEE Commun. Surv. Tutorials*, 22(1):746–789.
- Jonathan Hayase, Alisa Liu, Yejin Choi, Sewoong Oh, and Noah A. Smith. 2024. Data mixture inference: What do BPE tokenizers reveal about their training data? CoRR, abs/2407.16607.
- Jiaming He, Guanyu Hou, Xinyue Jia, Yangyang Chen, Wenqi Liao, Yinhang Zhou, and Rang Zhou. 2024. Data stealing attacks against large language models via backdooring. *Electronics*, 13(14):2858.
- Xuanli He, Lingjuan Lyu, Lichao Sun, and Qiongkai Xu. 2021. Model extraction and adversarial transferability, your BERT is vulnerable! In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 2006–2012. Association for Computational Linguistics.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531.
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2023. Unnatural instructions: Tuning language models with (almost) no human labor. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 14409–14428. Association for Computational Linguistics.
- Bo Hui, Haolin Yuan, Neil Gong, Philippe Burlina, and Yinzhi Cao. 2024. Pleak: Prompt leaking attacks against large language model applications. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, pages 3600–3614.
- IBM Security and Ponemon Institute. 2024. Cost of a data breach report 2024. https://www.ibm.com/reports/data-breach. Accessed: 2025-04-23.
- Daphne Ippolito, Nicholas Carlini, Katherine Lee, Milad Nasr, and Yun William Yu. 2023. Reverse-engineering decoding strategies given blackbox access to a language generation system. In *Proceedings of the 16th International Natural Language Generation Conference*, pages 396–406.
- Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot. 2020. High accuracy and high fidelity extraction of neural networks. In 29th USENIX Security Symposium, USENIX Security 2020, August 12-14, 2020, pages 1345–1362. USENIX Association.
- Mika Juuti, Sebastian Szyller, Samuel Marchal, and N. Asokan. 2019. PRADA: protecting against DNN

- model stealing attacks. In *IEEE European Symposium on Security and Privacy, EuroS&P 2019, Stockholm, Sweden, June 17-19, 2019*, pages 512–527. IEEE.
- Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter Kairouz, Christopher A. Choquette-Choo, and Zheng Xu. 2024. User inference attacks on large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 18238–18265. Association for Computational Linguistics.
- Manish Kesarwani, Bhaskar Mukhoty, Vijay Arya, and Sameep Mehta. 2018. Model extraction warning in mlaas paradigm. In *Proceedings of the 34th Annual Computer Security Applications Conference, ACSAC 2018, San Juan, PR, USA, December 03-07, 2018*, pages 371–380. ACM.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 17061–17084.
- Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A Hale. 2024. The benefits, risks and bounds of personalizing the alignment of large language models to individuals. *Nature Machine Intelligence*, 6(4):383–392.
- Kalpesh Krishna, Gaurav Singh Tomar, Ankur P. Parikh, Nicolas Papernot, and Mohit Iyyer. 2020. Thieves on sesame street! model extraction of bert-based apis. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.
- Andrew K. Lampinen, Ishita Dasgupta, Stephanie C. Y. Chan, Kory W. Mathewson, Michael Henry Tessler, Antonia Creswell, James L. McClelland, Jane Wang, and Felix Hill. 2022. Can language models learn from explanations in context? In Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022, pages 537–563. Association for Computational Linguistics.
- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. Textbugger: Generating adversarial text against real-world applications. In 26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019. The Internet Society.
- Marvin Li, Jason Wang, Jeffrey G. Wang, and Seth Neel. 2023. Mope: Model perturbation based privacy attacks on language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 13647–13660. Association for Computational Linguistics.

- Yiming Li, Linghui Zhu, Xiaojun Jia, Yong Jiang, Shu-Tao Xia, and Xiaochun Cao. 2022. Defending against model stealing via verifying embedded external features. In *Proceedings of the AAAI conference on* artificial intelligence, volume 36, pages 1464–1472.
- Zongjie Li, Chaozheng Wang, Pingchuan Ma, Chaowei Liu, Shuai Wang, Daoyuan Wu, Cuiyun Gao, and Yang Liu. 2024. On extracting specialized code abilities from large language models: A feasibility study. In Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024, pages 74:1–74:13.
- Zi Liang, Qingqing Ye, Yanyun Wang, Sen Zhang, Yaxin Xiao, Ronghua Li, Jianliang Xu, and Haibo Hu. 2024. Alignment-aware model extraction attacks on large language models. *CoRR*, abs/2409.02718.
- Allen Liu and Ankur Moitra. 2024. Model stealing for any low-rank language model. *Preprint*, arXiv:2411.07536.
- Daizong Liu, Mingyu Yang, Xiaoye Qu, Pan Zhou, Yu Cheng, and Wei Hu. 2024a. A survey of attacks on large vision-language models: Resources, advances, and future trends. *arXiv preprint arXiv:2407.07403*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.
- Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024b. Autodan: Generating stealthy jailbreak prompts on aligned large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. 2024c. Safety of multimodal large language models on images and text. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024*, pages 8151–8159. ijcai.org.
- Zhenhua Liu, Tong Zhu, Chuanyuan Tan, Bing Liu, Haonan Lu, and Wenliang Chen. 2024d. Probing language models for pre-training data detection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1576–1587. Association for Computational Linguistics.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards deep learning models resistant to adversarial attacks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 May 3, 2018, Conference Track Proceedings. OpenReview.net.
- Pratyush Maini, Hengrui Jia, Nicolas Papernot, and Adam Dziedzic. 2024a. LLM dataset inference: Did you train on my dataset? In *Advances in Neural*

- Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024.
- Pratyush Maini, Hengrui Jia, Nicolas Papernot, and Adam Dziedzic. 2024b. LLM dataset inference: Did you train on my dataset? In Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024.
- Chengzhi Mao, Ziyuan Zhong, Junfeng Yang, Carl Vondrick, and Baishakhi Ray. 2019. Metric learning for adversarial robustness. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 478–489.
- Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. 2023. Membership inference attacks against language models via neighbourhood comparison. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 11330–11343. Association for Computational Linguistics.
- Matthieu Meeus, Shubham Jain, Marek Rei, and Yves-Alexandre de Montjoye. 2024a. Did the neurons read your book? document-level membership inference for large language models. In 33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024. USENIX Association.
- Matthieu Meeus, Igor Shilov, Shubham Jain, Manuel Faysse, Marek Rei, and Yves-Alexandre de Montjoye. 2024b. Sok: Membership inference attacks on llms are rushing nowhere (and how to fix it). *arXiv* preprint arXiv:2406.17975.
- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date.
- Smitha Milli, Ludwig Schmidt, Anca D. Dragan, and Moritz Hardt. 2019. Model reconstruction from model explanations. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT** 2019, Atlanta, GA, USA, January 29-31, 2019, pages 1–9. ACM.
- Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. 2022. Quantifying privacy risks of masked language models using membership inference attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 8332–8347. Association for Computational Linguistics.
- John Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander Rush. 2023. Text embeddings reveal (almost) as much as text. In *Proceedings of the 2023*

- Conference on Empirical Methods in Natural Language Processing, pages 12448–12460. Association for Computational Linguistics.
- John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126.
- John X. Morris, Wenting Zhao, Justin T. Chiu, Vitaly Shmatikov, and Alexander M. Rush. 2024. Language model inversion. In *The Twelfth International Conference on Learning Representations, ICLR* 2024, *Vienna, Austria, May* 7-11, 2024. OpenReview.net.
- Hausi A. Müller, Jens H. Jahnke, Dennis B. Smith, Margaret-Anne D. Storey, Scott R. Tilley, and Kenny Wong. 2000. Reverse engineering: a roadmap. In 22nd International Conference on on Software Engineering, Future of Software Engineering Track, ICSE 2000, Limerick Ireland, June 4-11, 2000, pages 47–60. ACM.
- Ali Naseh, Kalpesh Krishna, Mohit Iyyer, and Amir Houmansadr. 2023. Stealing the decoding algorithms of language models. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, page 1835–1849.
- Milad Nasr, Javier Rando, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Florian Tramèr, and Katherine Lee. 2025. Scalable extraction of training data from aligned, production language models. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net.
- Myung Gyo Oh, Leo Hyun Park, Jaeuk Kim, Jaewoo Park, and Taekyoung Kwon. 2023. Membership inference attacks with token-level deduplication on korean language models. *IEEE Access*, 11:10207–10217.
- Seong Joon Oh, Bernt Schiele, and Mario Fritz. 2019. Towards reverse-engineering black-box neural networks. *Explainable AI: interpreting, explaining and visualizing deep learning*, pages 121–144.
- Daryna Oliynyk, Rudolf Mayer, and Andreas Rauber. 2023. I know what you trained last summer: A survey on stealing machine learning models and defences. *ACM Comput. Surv.*, 55(14s):324:1–324:41.
- OpenAI. 2020. Language models are few-shot learners. *CoRR*, abs/2005.14165.
- OpenAI. 2024. Introducing openai o1-preview: A new series of reasoning models for solving hard problems. available now.
- OpenAI. 2025. Rate limits openai api. https://platform.openai.com/docs/guides/rate-limits?utm_source=chatgpt.com.

- Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. 2019. Knockoff nets: Stealing functionality of blackbox models. In *IEEE Conference on Computer Vision* and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pages 4954–4963. Computer Vision Foundation / IEEE.
- Yonatan Oren, Nicole Meister, Niladri S. Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. 2024. Proving test set contamination in black-box language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Mustafa Özdayi, Charith Peris, Jack FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. 2023. Controlling the extraction of memorized data from large language models via prompt-tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 1512–1521. Association for Computational Linguistics.
- Soham Pal, Yash Gupta, Aditya Shukla, Aditya Kanade, Shirish K. Shevade, and Vinod Ganapathy. 2019. A framework for the extraction of deep neural networks by leveraging public data. *CoRR*, abs/1905.09165.
- Ashwinee Panda, Christopher A. Choquette-Choo, Zhengming Zhang, Yaoqing Yang, and Prateek Mittal. 2024. Teach llms to phish: Stealing private information from language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab Emirates, April 2-6, 2017*, pages 506–519. ACM.
- Rahil Parikh, Christophe Dupuy, and Rahul Gupta. 2022. Canary extraction in natural language understanding models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 552–560. Association for Computational Linguistics.
- Dario Pasquini, Evgenios M. Kornaropoulos, and Giuseppe Ateniese. 2024. Llmmap: Fingerprinting for large language models. *Preprint*, arXiv:2407.15847.
- Fábio Perez and Ian Ribeiro. 2022. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*.
- Ambrish Rawat, Stefan Schoepf, Giulio Zizzo, Giandomenico Cornacchia, Muhammad Zaid Hameed,

- Kieran Fraser, Erik Miehling, Beat Buesser, Elizabeth M. Daly, Mark Purcell, Prasanna Sattigeri, Pin-Yu Chen, and Kush R. Varshney. 2024. Attack atlas: A practitioner's perspective on challenges and pitfalls in red teaming genai. *CoRR*, abs/2409.15398.
- Scott E. Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. 2022. A generalist agent. *Trans. Mach. Learn. Res.*, 2022.
- Mark Russinovich and Ahmed Salem. 2024. Hey, that's my model! introducing chain & hash, an llm finger-printing technique. *Preprint*, arXiv:2407.10887.
- Amir Mahdi Sadeghzadeh, Amir Mohammad Sobhanian, Faezeh Dehghan, and Rasool Jalili. 2024. HODA: hardness-oriented detection of model extraction attacks. *IEEE Trans. Inf. Forensics Secur.*, 19:1429–1439.
- Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. 2019. Ml-leaks: Model and data independent membership inference attacks and defenses on machine learning models. In 26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019. The Internet Society.
- Zeyang Sha and Yang Zhang. 2024a. Prompt stealing attacks against large language models. *CoRR*, abs/2402.12959.
- Zeyang Sha and Yang Zhang. 2024b. Prompt stealing attacks against large language models. *arXiv* preprint *arXiv*:2402.12959.
- Adi Shamir, Isaac Andrés Canales Martinez, Anna Hambitzer, Jorge Chávez-Saab, Francisco Rodríguez-Henríquez, and Nitin Satpute. 2023. Polynomial time cryptanalytic extraction of neural network models. *CoRR*, abs/2310.08708.
- Virat Shejwalkar, Huseyin A Inan, Amir Houmansadr, and Robert Sim. 2021. Membership inference attacks against nlp classification models. In *NeurIPS 2021 Workshop Privacy in Machine Learning*.
- Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu, and Deyi Xiong. 2023. Large language model alignment: A survey. *CoRR*, abs/2309.15025.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024a. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-Review.net.

- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024b. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-Review.net.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pages 3–18. IEEE.
- Congzheng Song and Vitaly Shmatikov. 2019. Auditing data provenance in text-generation models. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 196–206.
- Tyler Sorensen and Heidy Khlaaf. 2024. Leftoverlocals: Listening to llm responses through leaked gpu local memory. *Preprint*, arXiv:2401.16603.
- Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15(1):1929–1958.
- Anshuman Suri. 2024. Reassessing emnlp 2024's best paper: Does divergence-based calibration for membership inference attacks hold up?
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014a. Intriguing properties of neural networks. In 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014b. Intriguing properties of neural networks. In 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Qwen Team. 2024. Qwq: Reflect deeply on the boundaries of the unknown.
- Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. 2018. Ensemble adversarial training: Attacks and defenses. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 May 3, 2018, Conference Track Proceedings. OpenReview.net.

- Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. 2016. Stealing machine learning models via prediction apis. In 25th USENIX Security Symposium, USENIX Security 16, Austin, TX, USA, August 10-12, 2016, pages 601–618.
- Tamás Várady, Ralph R. Martin, and Jordan Cox. 1997. Reverse engineering of geometric models an introduction. *Comput. Aided Des.*, 29(4):255–268.
- Eric Wallace, Mitchell Stern, and Dawn Song. 2020. Imitation attacks and defenses for black-box machine translation systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 5531–5546.
- Cheng Wang, Yiwei Wang, Bryan Hooi, Yujun Cai, Nanyun Peng, and Kai-Wei Chang. 2025. Con-recall: Detecting pre-training data in llms via contrastive decoding. In *Proceedings of the 31st International Conference on Computational Linguistics, COLING* 2025, Abu Dhabi, UAE, January 19-24, 2025, pages 1013–1026. Association for Computational Linguistics.
- Tianhao Wang and Florian Kerschbaum. 2021. Riga: Covert and robust white-box watermarking of deep neural networks. In *Proceedings of the Web Conference 2021*, pages 993–1004.
- Zhepeng Wang, Runxue Bao, Yawen Wu, Jackson Taylor, Cao Xiao, Feng Zheng, Weiwen Jiang, Shangqian Gao, and Yanfu Zhang. 2024. Unlocking memorization in large language models with dynamic soft prompting. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 9782–9796. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022.
- Junyi Wei, Yicheng Zhang, Zhe Zhou, Zhou Li, and Mohammad Abdullah Al Faruque. 2020. Leaky dnn: Stealing deep-learning model secret with gpu context-switching side-channel. In 2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), pages 125–137. IEEE.
- Guanlong Wu, Zheng Zhang, Jianyu Niu, Weili Wang, Yao Zhang, Ye Wu, and Yinqian Zhang. 2025. I know what you asked: Prompt leakage via kv-cache sharing in multi-tenant llm serving. *Network and Distributed System Security Symposium*.
- Roy Xie, Junlin Wang, Ruomin Huang, Minxing Zhang, Rong Ge, Jian Pei, Neil Gong, and Bhuwan Dhingra.

- 2024. Recall: Membership inference via relative conditional log-likelihoods. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 8671–8689. Association for Computational Linguistics.
- Cheng Xiong, Guorui Feng, Xinran Li, Xinpeng Zhang, and Chuan Qin. 2022. Neural network model protection with piracy identification and tampering localization capability. In *Proceedings of the 30th ACM International Conference on Multimedia*, page 2881–2889, New York, NY, USA.
- Jiashu Xu, Fei Wang, Mingyu Ma, Pang Wei Koh, Chaowei Xiao, and Muhao Chen. 2024. Instructional fingerprinting of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3277–3306.
- Haonan Yan, Xiaoguang Li, Hui Li, Jiamin Li, Wenhai Sun, and Fenghua Li. 2022a. Monitoring-based differential privacy mechanism against query floodingbased model extraction attack. *IEEE Trans. Dependable Secur. Comput.*, 19(4):2680–2694.
- Haonan Yan, Xiaoguang Li, Hui Li, Jiamin Li, Wenhai Sun, and Fenghua Li. 2022b. Monitoring-based differential privacy mechanism against query floodingbased model extraction attack. *IEEE Trans. Depend*able Secur. Comput., 19(4):2680–2694.
- Mengjia Yan, Christopher W Fletcher, and Josep Torrellas. 2020. Cache telepathy: Leveraging shared resource attacks to learn {DNN} architectures. In 29th USENIX Security Symposium (USENIX Security 20), pages 2003–2020.
- Kang Yang, Run Wang, and Lina Wang. 2022. Metafinger: Fingerprinting the deep neural networks with meta-training. In *IJCAI*, pages 776–782.
- Wentao Ye, Jiaqi Hu, Liyao Li, Haobo Wang, Gang Chen, and Junbo Zhao. 2024. Data contamination calibration for black-box llms. In *Findings of the Association for Computational Linguistics*, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pages 10845–10861. Association for Computational Linguistics.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. Privacy risk in machine learning: Analyzing the connection to overfitting. In 31st IEEE Computer Security Foundations Symposium, CSF 2018, Oxford, United Kingdom, July 9-12, 2018, pages 268–282. IEEE Computer Society.
- Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. 2023. Bag of tricks for training data extraction from language models. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 40306–40320. PMLR.

Santiago Zanella-Beguelin, Shruti Tople, Andrew Paverd, and Boris Köpf. 2021. Grey-box extraction of natural language models. In *Proceedings of the 38th International Conference on Machine Learning*, pages 12278–12286.

Boyi Zeng, Lizheng Wang, Yuncong Hu, Yi Xu, Chenghu Zhou, Xinbing Wang, Yu Yu, and Zhouhan Lin. 2024. Huref: HUman-REadable fingerprint for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*

Chaoning Zhang, Chenshuang Zhang, Sheng Zheng, Yu Qiao, Chenghao Li, Mengchun Zhang, Sumit Kumar Dam, Chu Myaet Thwal, Ye Lin Tun, Le Luang Huy, Dong Uk Kim, Sung-Ho Bae, Lik-Hang Lee, Yang Yang, Heng Tao Shen, In So Kweon, and Choong Seon Hong. 2023. A complete survey on generative AI (AIGC): is chatgpt from GPT-4 to GPT-5 all you need? *CoRR*, abs/2303.11717.

Collin Zhang, John X. Morris, and Vitaly Shmatikov. 2024a. Extracting prompts by inverting LLM outputs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 14753–14777. Association for Computational Linguistics.

Weichao Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024b. Pretraining data detection for large language models: A divergence-based calibration method. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 5263–5274. Association for Computational Linguistics.

Weichao Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024c. Pretraining data detection for large language models: A divergence-based calibration method. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 5263–5274. Association for Computational Linguistics.

Yiming Zhang, Nicholas Carlini, and Daphne Ippolito. 2024d. Effective prompt extraction from language models. In *First Conference on Language Modeling*.

Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. 2022. Decoupled knowledge distillation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR* 2022, New Orleans, LA, USA, June 18-24, 2022, pages 11943–11952. IEEE.

Xinyao Zheng, Husheng Han, Shangyi Shi, Qiyan Fang, Zidong Du, Xing Hu, and Qi Guo. 2024. Inputsnatch: Stealing input in llm services via timing side-channel attacks. *Preprint*, arXiv:2411.18191.

Yuankun Zhu, Yueqiang Cheng, Husheng Zhou, and Yantao Lu. 2021. Hermes attack: Steal {DNN} models with lossless inference accuracy. In 30th USENIX Security Symposium (USENIX Security 21).

A Knowledge Distillation

We have introduced function imitation in section 3.2.3, to some extent, knowledge distillation (Hinton et al., 2015) can also be considered as a form of function imitation. Knowledge Distillation aims to transfer knowledge from a large teacher model to a small student model. By encouraging the student model to approximate the behavior of the teacher model, the student is able to achieve functional imitation with minimal loss in quality, while achieving higher inference efficiency (Zhao et al., 2022; Gu et al., 2024).

However, most knowledge distillation methods often assume white-box access to the teacher model and have a certain understanding of training data. Therefore, its application on model reverse is limited, and it is not considered as a primary attack method.

B Function Imitation for DNNs

For neural network models, the extraction of them is relatively easier compared to transformer models. As a result, attackers typically do not achieve function imitation by training imitation model with input-output pairs, but instead analyze parameters directly and then reconstruct the model. Milli et al. (2019); Jagielski et al. (2020) covers some common strategies for DNNs imitation, through multiple queries and algebraic methods, attackers can estimate the number of layers, the activation functions used, and the overall structure of the model. Pal et al. (2019); Shamir et al. (2023), through using activate learning to effectively generate queries, select the most informative samples for better reconstruction.

C Reverse Engineering of Multimodel Large Language Models

The emergence of Multimodal Large Language Models (MLLMs) (Liu et al., 2023; Achiam et al., 2023; Team et al., 2023) has introduced both new opportunities and unique challenges in the context of reverse engineering. Unlike traditional language models, MLLMs process not only textual data but also other modalities such as images, audio, and video, creating additional attack surfaces. Similar to other attacks (Liu et al., 2024c,a), these expanded interfaces are anticipated to heighten the models' susceptibility to reverse engineering attempts. For instance, the integration of visual inputs, such as images, presents new challenges, in-

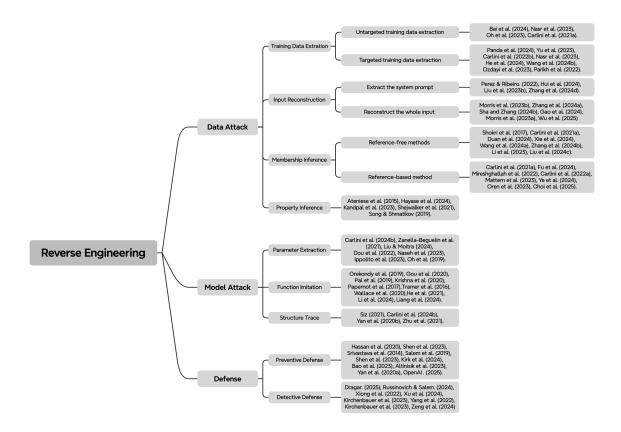


Figure 5: A taxonomy of the paper

cluding adversarial visual perturbations (Szegedy et al., 2014b; Madry et al., 2018), which can be more more dangerous and difficult to mitigate (Carlini et al., 2023b) compared to adversarial textual perturbations (Morris et al., 2020; Li et al., 2019).

Therefore, future research could explore interesting topics such as: • Benchmarking vulnerabilities of MLLMs to reverse engineering. • Developing strategies of multimodal reverse engineering. • Designing robust protective mechanisms.

D Extra Experiment for Latest Model

First, we survey the targets of the latest and most representative reverse-engineering techniques, as summarized in Table 1. The data reveal that most attacks focus on open-source models, while among commercial offerings, current efforts concentrate predominantly on GPT-3.5 Turbo. This disparity arises partly from the ease of evaluating attack efficacy on open-source platforms and partly from the more comprehensive defenses employed by commercial providers. Accordingly, a systematic assessment of these methods' performance on state-of-the-art models is both warranted and valuable

for guiding future research.

Therefore, we compiled a collection of representative reverse engineering studies with actively maintained codebases and evaluated their methods on GPT-40. We note that the membership inference attack experiments are detailed in the following section.

For the training data extraction phase, we selected three methods from (Carlini et al., 2021a; Özdayi et al., 2023; Bai et al., 2024). Although evaluating the success of data extraction attacks is inherently challenging, our experiments show that these techniques failed to recover any meaningful information, yielding virtually no outputs resembling the original training data.

For the prompt extraction and property inference phase, we evaluated four methods from (Perez and Ribeiro, 2022; Hui et al., 2024; Zhang et al., 2024d; Hayase et al., 2024). Our results show that, relative to training data recovery, these prompt extraction techniques achieve substantially higher success rates. However, it is worth noting that prompt defenses have evolved just as quickly: OpenAI is progressively deploying countermeasures against prompts that exhibit high extraction success rates.

Table 2: Model targets of some newest attack

| Attack Type | Method | GPT-2 | Falcon | Pythia | Llama | Llama-2 | Llama-3 | Mistral | GPT-3.5-turbo | GPT-40 |
|--------------------|--|----------|----------|----------|----------|---------|---------|----------|---------------|--------|
| Training Data | Carlini et al. (2021a) Nasr et al. (2025) Bai et al. (2024) Panda et al. (2024) | V V | V V | V | V | V | V | V | V V | |
| Prompt Extract | Hui et al. (2024) Sha and Zhang (2024a) | | V | | V | ✓ | | | V | |
| Property Inference | Hayase et al. (2024) | V | | | V | | V | V | ✓ | |
| MIA | Maini et al. (2024a) | | | / | | | | | | |
| Model parameter | Carlini et al. (2024) | V | | V | / | | | | ✓ | |
| Model function | Li et al. (2024) | | | | | | | | ✓ | |

Table 3: Evaluation of Existing Attack Methods

| Attack Type | Method | dataset | Effectiveness | Prerequisites | Query Count | Leakage Quality |
|--------------------|--------------------------|-------------------------|---------------|---------------|-------------|-----------------|
| | Carlini et al. (2021a) | | X | | | |
| Training Data | Özdayi et al. (2023) | | × | | | |
| | Bai et al. (2024) | | X | | | |
| | Perez and Ribeiro (2022) | | X | | | |
| Prompt Extract | Hui et al. (2024) | | ✓ | low | low | medium |
| | Zhang et al. (2024d) | awesome-chatgpt-prompts | ✓ | low | low | high |
| Property Inference | Hayase et al. (2024) | Oscar | ✓ | high | low | high |
| Model parameter | Carlini et al. (2024) | | X | | | |
| Model function | Li et al. (2024) | | ✓ | low | high | high |

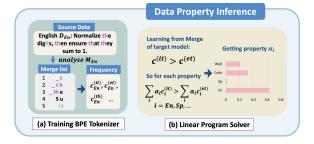


Figure 6: The illustration of the data property inference attack, where most commercial models publicly release their merge.txt file and the source data comprise diverse datasets hosted on Hugging Face.

Table 4: Experiment on Effective Prompt Extraction from Models

| Dataset | awesome | sharegpt | unnatural |
|---------|---------|----------|-----------|
| exact | 54.1 | 48.1 | 68.2 |
| approx | 81.3 | 79.4 | 74.8 |

Due to the high computational cost and the absence of publicly available code in most model level attack studies, we selected two representative methods for our experiments (Carlini et al., 2024; Li et al., 2024). We note that, because few

Table 5: Property Inference of GPT-40

| Category | GPT-40 | LLAMA 3 | Claude |
|----------|--------|---------|--------|
| Web | 20.5 | 12.7 | 25.6 |
| Code | 32.8 | 30.3 | 25.8 |
| Book | 7.4 | 8.5 | 12.8 |
| French | 2.9 | 1.8 | 3.1 |

security papers provide complete implementations, we effectively executed every technique with sufficient supporting code or detailed descriptions. As demonstrated above, many of these approaches have since been mitigated by (i) more restrictive access policies, (ii) accelerated vulnerability patching, and (iii) increasingly robust defense mechanisms, rendering them largely ineffective against today's state of the art models. Nonetheless, their foundational insights remain valuable: data reconstruction and functionality extraction can be further refined through additional experimentation, while full model reconstruction continues to pose an open research challenge, one that will require substantial future investment and resource intensive efforts.

E Current Problems in MIA

Although membership inference attacks were first proposed by Shokri et al. (2017) and validated on classifiers and fine-tuned models, recent papers (Duan et al., 2024; Meeus et al., 2024b; Maini et al., 2024b) and blog posts (Suri, 2024) have shown their underwhelming performance on pretrained large-scale models. Motivated by these findings, we conducted some simple experiments on the Pythia-1.4B to intuitively expose potential shortcomings in current MIA methodologies, datasets, and benchmarking practices, as we show in Table 6, compared to the randomly partitioned Wikipedia dataset, WikiMIA exhibits pronounced distributional drift. In drift-free datasets, the four MIA techniques—loss-based (Yeom et al., 2018), reference-based (Carlini et al., 2021b), Mink (Shi et al., 2024b), and zlib (Carlini et al., 2021c)—achieve near-random membership inference; however, their efficacy notably increases on GitHub data, revealing dataset-specific biases—for example, zlib performs best on GitHub but falls short of Ref on Wikipedia. This motivated us to systematically summarize the existing challenges in the MIA field:

Table 6: Traditional MIA method on LLM

| Category | Loss | Ref | min-k | zlib |
|-----------|-------|-------|-------|-------|
| Wikimia | 0.534 | 0.607 | 0.685 | 0.674 |
| Wikipedia | 0.516 | 0.571 | 0.514 | 0.524 |
| Github | 0.654 | 0.594 | 0.643 | 0.671 |

(i) Improper membership splitting. Instead of random sampling, some studies construct member and non-member sets post hoc-after model training—using non-random criteria such as corpus origin, timestamps, or labels. This practice introduces severe distributional drift and semantic cue leakage, causing attacks to exploit differences in writing style or token frequencies rather than true membership signals. For example, a 2023 corpus contains time-sensitive tokens like "COVID-19" or "Chat-GPT" that are absent in a 2020 dataset, allowing MIAs to distinguish samples based solely on their relative occurrence frequencies. Duan et al. (2024) conducted more detailed experiments and showed that fuzzy leakage can occur even when there is no exact overlap between member and non-member samples. They argue that semantic and syntactic similarity measures should be incorporated into the design of more robust evaluation frameworks

and benchmarks. Meeus et al. (2024b) also point out that certain datasets—such as WikiMIA, arXiv, Books, and Stack—may inherently exhibit distributional drift, which undermines the reliability of results derived from them.

- (ii) Excessive pretraining scale. Large language models are trained for just one epoch over massive corpora, which dilutes their retention of individual samples. As a result, many attack assumptions that hold for classifiers break down on LLMs, that's why loss-based inference methods perform at near-chance levels in MIA evaluations against large pre-trained models. To address the scale and industrial requirements of modern LLMs, Maini et al. (2024b) extend the membership inference paradigm to the dataset level and introduce a novel detection framework—dataset inference—which employs a composite indicator function to determine whether a given dataset was used in the model's pretraining.
- (iii) Lack of standardized benchmarks and protocols. Studies often employ disparate models and evaluation suites without common control experiments, and attack performance varies across domains. This inconsistency makes it difficult to quantitatively compare the effectiveness of different MIA methods.

F Frequently Chosen Benchmarks for Data Recovery Engine

We have collected frequently used metrics in Table 7 and datasets in Table 8.

G Prompt Extraction Examples

Here, we present successful cases of prompt extraction against several state-of-the-art commercial models, as shown in Table ??. Furthermore, our experiments reveal that reasoning oriented models exhibit markedly greater resilience to prompt extraction attacks: most prompts that succeed against GPT-40 are ineffective when applied to these reasoning models.

H Comparison with Existing Surveys

For the multiple surveys mentioned in the Introduction, we provide here a detailed comparison to highlight how they differ from our work (as shown in 10).

¹https://github.com/google-research/ lm-extraction-benchmark

²https://github.com/f/awesome-chatgpt-prompts

³https://github.com/sahil280114/codealpaca

Table 7: Frequently evaluated metrics of data recovery engine. MIA, TDEA and IRA stand for membership inference attack, training data extraction attack and input reconstruction attack.

| Attack Type | Metric Name | Explanation |
|--------------|-----------------------------------|--|
| MIA MIA | AUROC TPR@5% FPR | Area under the ROC curve. true positive rate at 5% false positive rate. |
| TDEA TDEA | Speed Verbatim Extraction Rate | The amount of time required for the attack to execute. The proportion of generated suffixes that exactly match the original text. |
| TDEA | Paraphrase Extraction Rate | The proportion of generated suffixes that are rephrased versions of the original text. |
| IRA | BLEU | N-gram similarity between the original and reconstructed texts. |
| IRA | Exact Match | The multi-class F1 score comparing the set of predicted tokens to the set of true tokens. |
| IRA | Token-level F1 Score | The percentage of reconstructed outputs that exactly match the ground truth. |
| IRA | Semantic Similarity | The cosine similarity between the output of the text embedding models. |

Table 8: Frequently evaluated datasets of data recovery engine. MIA, TDEA and IRA stand for membership inference attack, training data extraction attack and input reconstruction attack.

| Attack Type | Dataset Name | Comment |
|-------------|--|--|
| MIA | WikiMIA (Shi et al., 2024a) | split sentence-level members/non-members by date. |
| MIA | StackMIAsub (Ye et al., 2024) | split sentence-level members/non-members by date. |
| MIA | MIMIR (Duan et al., 2024) | split sentence-level members/non-members by original training/test set. |
| MIA | ArXiv - Document (Meeus et al., 2024a) | split document-level members/non-members by date. |
| TDEA | LM Extraction Benchmark ¹ | prefixes if 50-token length are given to extract the suffixes based on The Pile dataset. |
| IRA | Alpaca Code Generation ² | Code prompts from Alpaca. |
| IRA | Awesome-ChatGPT-Prompts ³ | Detailed prompts designed to adapt the LLM to a specific role. |
| IRA | Unnatural Instructions (Honovich et al., 2023) | A large, diverse set of instructions, collected with minimal human effort. |

Table 10: Comparison with Existing Surveys in topic and scope.

| Ref | Year | Topic and Motivation | Covered Technical Approaches |
|-----------------------|------|--|--|
| Esmradi et al. (2023) | 2023 | Analyze different types of attacks targeting LLMs. | Prompt Injection, Privacy Leakage, Model Theft, Data Reconstruction, Data Poisoning, Model Hijack- ing, Membership Inference |
| Das et al. (2025) | 2024 | Organize privacy and security challenges in LLMs. | Prompt Hacking, Adversarial Attacks, Gradient Leakage, Membership Inference, PII Leakage |
| Chen et al. (2025) | 2025 | Discuss privacy risks and protection methods. | Backdoor Attacks, Model Inversion, Model Stealing, Data Stealing, Training Data Extraction, Membership & Attribute Inference |
| Ours | 2025 | Recover/infer protected information from LLMs. | Training Data Extraction, Input Reconstruction, Membership & Property Inference, Model Parameter Extraction, Function Imitation, Structure Tracing |