



Refinement Provenance Inference: Detecting LLM-Refined Training Prompts from Model Behavior

Anonymous ACL submission

Abstract

Instruction tuning increasingly relies on LLM-based prompt refinement, where prompts in the training corpus are selectively rewritten by an external refiner to improve clarity and instruction alignment. This motivates an instance-level audit problem: for a fine-tuned model and a training prompt–response pair, can we infer whether the model was trained on the original prompt or its LLM-refined version within a mixed corpus? This matters for dataset governance and dispute resolution when training data are contested. However, it is non-trivial in practice: refined and raw instances are interleaved in the training corpus with unknown, source-dependent mixture ratios, making it harder to develop provenance methods that generalize across models and training setups. In this paper, we formalize this audit task as Refinement Provenance Inference (RPI) and show that prompt refinement yields stable, detectable shifts in teacher-forced token distributions, even when semantic differences are not obvious. Building on this phenomenon, we propose RePro, a logit-based provenance framework that fuses teacher-forced likelihood features with logit-ranking signals. During training, RePro learns a transferable representation via shadow fine-tuning, and uses a lightweight linear head to infer provenance on unseen victims without training-data access. Empirically, RePro consistently attains strong performance and transfers well across refiners, suggesting that it exploits refiner-agnostic distribution shifts rather than rewrite-style artifacts.

1 Introduction

Large language models have rapidly evolved into general-purpose systems that power a wide range of applications, from reasoning and code generation to dialogue and tool use (Achiam et al., 2023; Dubey et al., 2024; Team et al., 2023; Roziere et al., 2023). As capabilities have improved, fine-tuning and instruction tuning have become standard prac-

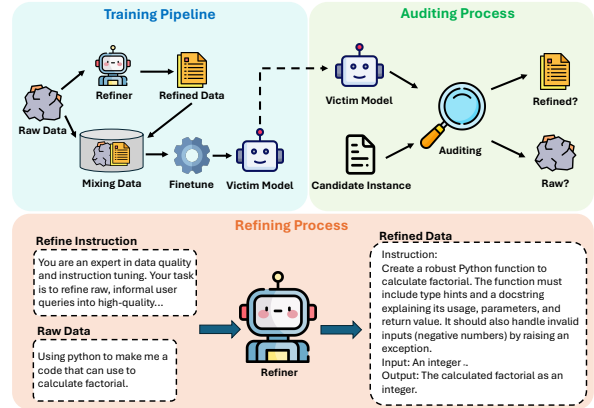


Figure 1: The process of Refinement Provenance Inference (RPI) problem.

tics for adapting these models to specific domains and interaction styles, which in turn has placed increasing emphasis on the construction and curation of high-quality training prompts (Yin et al., 2025a,b; Ouyang et al., 2022; Hu et al., 2022; Li et al., 2023a). In many modern pipelines, raw prompts, the original collected instructions, are rewritten by a refiner model to standardize phrasing, reduce ambiguity, and align with instruction-following conventions (Xu et al., 2024a; Mukherjee et al., 2023; Yan et al., 2025; Liu et al., 2024). This widely used refinement step raises a provenance question for auditing: given a fine-tuned model and a candidate instance (x_j, y_j) , **can we infer whether the model was tuned on its raw version or on its refiner-rewritten counterpart?**

Answering this question matters for both transparency and risk assessment in model development (Longpre et al., 2023; Mitchell et al., 2019). From an auditing perspective, refinement can materially change the distribution of training prompts, and practitioners may wish to verify whether a deployed model was trained under a declared data pipeline or whether an undisclosed refiner was used (Mökander et al., 2024; Dziedzic et al., 2022). From a security and privacy perspective, the refinement step may also act as a distinctive trans-

071	formation that leaks information about the training	Our contributions are as follows:	123
072	process itself, potentially revealing aspects of an		
073	organization’s data preparation workflow (Carlini	• New provenance task. We introduce Re-	124
074	et al., 2022b; Nasr et al., 2023; Li et al., 2025a).	finement Provenance Inference (RPI), which	125
075	We refer to this auditing problem as Refinement	asks, for a candidate instance (x, y) and a	126
076	Provenance Inference (RPI). Figure 1 shows the	fine-tuned model, whether the model’s fine-	127
077	process of RPI. Notably, this is a data-level prove-	tuning data used the raw prompt or its refiner-	128
078	nance problem: the victim may have been fine-	rewritten version for that instance, and we	129
079	tuned on a mixture of refined and raw prompts, and	frame it as an actionable auditing problem for	130
080	the goal is to localize which training instances were	modern fine-tuning pipelines.	131
081	refined. A natural hypothesis is that refinement	• Provenance framework. We propose Re-	132
082	primarily changes surface form, and that prove-	Pro, a logit-based provenance framework that	133
083	nance evidence would therefore be tied to the spe-	extracts complementary teacher-forced logit	134
084	cific refiner and its rewriting style. However, we	cues and learns a transferable embedding via	135
085	argue that training on refined prompts induces	shadow fine-tuning and supervised contrastive	136
086	distribution-level preference shifts that persist	learning, enabling inference on victim models	137
087	beyond surface realizations. Concretely, refine-	using a lightweight linear classifier.	138
088	ment tends to make prompts more canonical and	• Evidence for detectable and transferable	139
089	better aligned with instruction-following conven-	traces. We provide a comprehensive empir-	140
090	tions, which biases the gradients observed during	ical study across tasks, victim families, and	141
091	fine-tuning and can alter the victim model’s token-	refinement operators, including cross-refiner	142
092	level preferences under teacher forcing (Shumailov	transfer, feature and training ablations, and	143
093	et al., 2023; Zhou et al., 2023; Santurkar et al.,	sensitivity analyses that characterize when re-	144
094	2023; Gudibande et al., 2023). These shifts are	finement traces are detectable and which com-	145
095	not always obvious from generated text, but they	ponents drive performance.	146
096	can be measured from the teacher-forced token dis-		
097	tributions as changes in likelihood patterns, rank-	2 Related Work	147
098	ing behavior among top candidates, and logit mar-		
099	gins (Shi et al., 2023; Hans et al., 2024; Gonen	2.1 Training-Data Auditing	148
100	et al., 2023). The central challenge is to extract sig-		
101	nals that are robust to variation in victim families	A long line of work shows that training induces sys-	149
102	and refinement operators, and to do so in a way that	tematic changes in a model’s confidence landscape	150
103	transfers across models rather than overfitting to a	that can be exploited for auditing (Shokri et al.,	151
104	particular refiner or data distribution.	2017; Yeom et al., 2018; Song and Mittal, 2021;	152
105	To address this challenge, we propose RePro ,	Salem et al., 2018). In membership inference, at-	153
106	a logit-based framework for refinement prove-	tackers distinguish seen versus unseen examples	154
107	nance inference that learns transferable signals in a	using statistics such as loss, entropy, or margin,	155
108	shadow training setup. We first compute a compact	and stronger variants rely on shadow-model trans-	156
109	feature vector from teacher-forced logits, captur-	fer, calibration features, or query-efficient prob-	157
110	ing complementary evidence such as token-level	ing (Carlini et al., 2022a; Duan et al., 2024; Ko	158
111	negative log-likelihood statistics, ranking patterns	et al., 2023; Li et al., 2025b). Beyond membership,	159
112	among top candidates, and margin features derived	property inference predicts whether the training	160
113	from logit gaps, with an uplift. We then train an em-	set contains examples with a particular attribute	161
114	bedding encoder via supervised contrastive learn-	by aggregating output statistics, highlighting that	162
115	ing on shadow models fine-tuned from the same	model outputs can leak training-time signals even	163
116	base initialization, encouraging embeddings with	when the attribute is not directly observable from	164
117	the same provenance label to cluster while separ-	the generated text (Atenièse et al., 2015; Kandpal	165
118	ating embeddings with different labels. Finally, we	et al., 2024; Ganju et al., 2018; Mahloujifar et al.,	166
119	fit a lightweight linear classifier on top of the frozen	2022). For language models, studies on memo-	167
120	embeddings and transfer the resulting attacker to	rization and data extraction further support that	168
121	victim models to produce refined-versus-raw prove-	token-level likelihood patterns can encode training-	169
122	nance scores at inference time.	time regularities (Carlini et al., 2022b; Shi et al.,	170
		2023). Our work follows this general paradigm	171

but targets a different training attribute, namely whether a fine-tuning instance was presented in an LLM-refined form rather than its raw form, which motivates logit-centric signals and transfer-based attackers instead of text-only evidence.

2.2 Data Refinement and Detection

LLM-driven rewriting is now widely used in large-scale data curation, particularly for instruction-tuning where prompts are standardized, clarified, and aligned to target interaction styles (Xu et al., 2024b; Ding et al., 2023; Li et al., 2024). Prior work studies refinement operators and policies, showing that automated rewriting can shift both surface form and latent preferences, yielding refined distributions that systematically differ from raw data (Lee et al., 2023; Sun et al., 2023; Li et al., 2023b). While refinement is typically treated as a quality-improving preprocessing step, its downstream footprint as an auditable training attribute has received less attention (Golchin and Surdeanu, 2023; Zhang et al., 2024; Lyu and Yin, 2024). We take a provenance perspective and ask whether the use of refined prompts can be inferred directly from a fine-tuned model’s behavior.

A related line of research aims to detect machine-generated or machine-transformed text via likelihood artifacts, perturb-and-score stability tests, stylistometric signals, and watermarking (Mitchell et al., 2023; Yang et al., 2023; Su et al., 2023; Kirchenbauer et al., 2023; Kuditipudi et al., 2023). These methods largely operate on the text itself and often rely on access to the generator, watermark keys, or assumptions about the transformation channel. Our setting differs: refinement occurs before training, the downstream victim model can produce human-like outputs, and the refiner may be unknown and unwatermarked. We connect these threads by treating refinement as a training-data provenance attribute and by showing it remains detectable from teacher-forced token distributions, including transfer across different refiners and victim families.

3 Refinement Provenance Inference

3.1 Problem Definition

Modern fine-tuning pipelines often refine training prompts using an external LLM to improve clarity and consistency. Such refinement can induce systematic shifts in the effective training distribution, which may be reflected in the token-level predictive behavior of the fine-tuned model. We study refine-

ment provenance inference at the instance level: within a single fine-tuned victim, different training instances may use different prompt variants, and the goal is to infer, for each instance, whether the prompt used during fine-tuning was raw or LLM-refined. We emphasize that refinement is applied only to the input prompt, while the reference output remains unchanged.

We index semantic instances by i , where each instance corresponds to a unique underlying task with a raw prompt x_i^{raw} and a reference output y_i . A refinement operator $R(\cdot)$ maps the raw prompt to a refined prompt:

$$x_i^{\text{ref}} = R(x_i^{\text{raw}}). \quad (1)$$

A fine-tuning dataset is constructed by selecting, for each instance i , either the raw or refined prompt with an i.i.d. latent indicator $z_i \sim \text{Bernoulli}(\rho)$:

$$x_i^{\text{tr}} = x_i^{z_i} = \begin{cases} x_i^{\text{ref}}, & z_i = 1, \\ x_i^{\text{raw}}, & z_i = 0, \end{cases} \quad z_i \in \{0, 1\}, \quad (2)$$

where z_i is the refinement provenance label (1 for refined, 0 for raw). Let M_0 denote a base language model and M_a denote the victim obtained by supervised fine-tuning (SFT) on the mixture:

$$M_a \leftarrow \text{SFT}(M_0; \{(x_i^{\text{tr}}, y_i)\}_{i \in \mathcal{I}_a}), \quad (3)$$

where \mathcal{I}_a is the set of semantic instances used to fine-tune the victim.

Auditing task. For an instance $i \in \mathcal{I}_a$ (membership known), the auditor is given the victim M_a and an evaluation pair (\tilde{x}_i, y_i) for teacher forcing, and aims to infer the training-time provenance label z_i :

$$s_i = g(\phi(M_a; \tilde{x}_i, y_i)) \in [0, 1], \quad \hat{z}_i = \mathbf{I}[s_i \geq \tau], \quad (4)$$

where $\phi(\cdot)$ extracts features from the victim’s token-level predictive behavior on y_i conditioned on \tilde{x}_i , $g(\cdot)$ outputs a classification score, and τ is a threshold. In our main setting, $\phi(\cdot)$ is computed from teacher-forced log-probabilities and top- k logit statistics.

3.2 Access Assumptions

We assume an auditor has (i) query access to a fine-tuned victim model M_a ; (ii) an evaluation set of instances with reference outputs; and (iii) access to the underlying base model M_0 to construct shadow fine-tuned models for learning transferable decision rules. Given a candidate pair (x_j, y_j) , the auditor

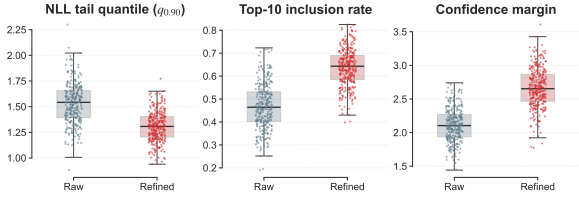


Figure 2: Feature diagnostics for teacher-forced logit.

performs teacher forcing on y_j to obtain token-level log probabilities and compute NLL-based statistics. Our main setting further assumes access to top- k logits. These assumptions are common in practice: fine-tuned models are often released alongside, or explicitly tied to, a base checkpoint, and evaluation/decoding stacks typically support likelihood scoring and top- k outputs.

4 Methodology

Our goal is to infer whether a fine-tuned model M_a was trained on a raw or an LLM-refined version of an instance. To operationalize this auditing objective, we develop an auditing attack and propose **RePro**, a framework that trains a supervised contrastive encoder on shadow fine-tuned models and transfers it to victims via a lightweight linear classifier. Figure 3 provides an overview of the full pipeline. In Stage 1, we construct a labeled shadow mixture of raw and refined instances, fine-tune a shadow model M_c from the same base model M_0 , and extract logit-derived feature vectors that summarize teacher-forced behavior through complementary signals such as NLL statistics, Top- K ranking patterns, logit margins, and optional uplift features. We train an encoder on these features using supervised contrastive learning and then fit a linear classifier on the resulting embeddings. In Stage 2, we apply the same feature extraction and encoder to a victim model M_a and use the transferred classifier to output a refined-versus-raw provenance score for each candidate instance.

4.1 Teacher-Forced Logit Features

Given a model M and an instance (x_i, y_i) , we compute token-level log-probabilities under teacher forcing:

$$\ell_{i,t}^{(M)} = \log p_M(y_{i,t} | x_i, y_{i,<t}), \quad t = 1, \dots, |y_i|. \quad (5)$$

Let $s_{i,t}^{(M)} \in \mathbb{R}^{|\mathcal{V}|}$ denote the pre-softmax logit vector at step t . From $\{\ell_{i,t}^{(M)}\}_{t=1}^{|y_i|}$ and the corresponding logits, we derive a fixed-dimensional feature vector that summarizes (i) likelihood-based fit and

tail hardness, and (ii) logit-based ranking behavior and local distribution sharpness.

Normalized negative log-likelihood (NLL).

$$\text{NLL}_M(i) = -\frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \ell_{i,t}^{(M)}. \quad (6)$$

To capture “hard-token” tails that are obscured by averaging, we additionally compute selected quantiles of tokenwise NLL values $\{-\ell_{i,t}^{(M)}\}$.

Top- k inclusion. Let $\text{TopK}_t^{(M)}$ denote the set of top- k tokens under M at step t (ranked by logits in $s_{i,t}^{(M)}$). We define

$$\text{TopK}_M(i) = \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \mathbb{I}[y_{i,t} \in \text{TopK}_t^{(M)}], \quad (7)$$

and use $k \in \{1, 5, 10\}$ in our experiments. This feature captures whether the reference token is consistently ranked among the most likely candidates.

Confidence margin. Let $s_{t,(1)}^{(M)}$ and $s_{t,(2)}^{(M)}$ be the largest and second-largest logit values in $s_{i,t}^{(M)}$, respectively. We compute the average margin

$$\text{Gap}_M(i) = \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} (s_{t,(1)}^{(M)} - s_{t,(2)}^{(M)}). \quad (8)$$

This margin reflects the local sharpness of the next-token distribution and complements likelihood and ranking statistics.

Uplift. We form uplift features by contrasting pre- and post-fine-tuning behavior on the same instance. For any statistic $S(\cdot) \in \{\text{NLL}, \text{TopK}, \text{Gap}\}$, we define

$$\Delta S(i) = S_{M_0}(i) - S_M(i), \quad (9)$$

where M is the fine-tuned model of interest (victim M_a or shadow M_c). Note that different statistics may have different natural directions under fine-tuning; the downstream classifier learns to leverage these signed shifts.

We aggregate the above statistics into a feature vector

$$\phi(M; x_i, y_i) \in \mathbb{R}^d, \quad (10)$$

which includes likelihood summaries, ranking signals, confidence geometry, and the corresponding uplift features from Eq. (9). Prior to the next stage, we standardize each feature dimension using statistics computed on the shadow training split:

$$\tilde{\phi}_j = \frac{\phi_j - \mu_j}{\sigma_j + \epsilon}, \quad (11)$$

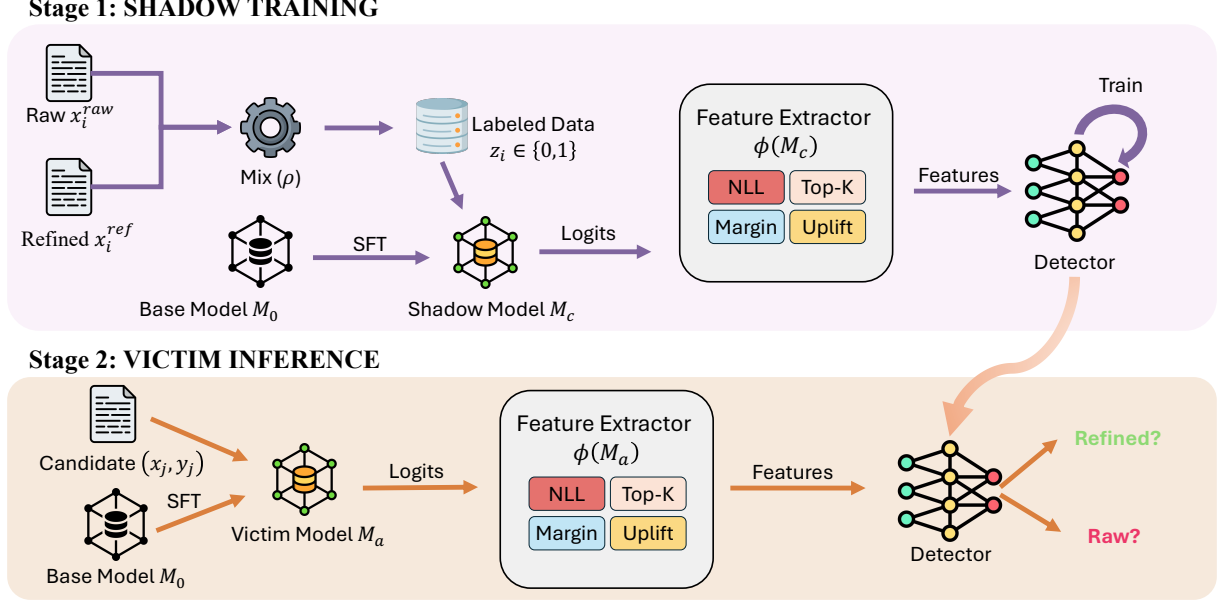


Figure 3: Overview of RePro. Stage 1 trains a supervised-contrastive encoder on logit-derived features from a shadow fine-tuned model. Stage 2 transfers the encoder and a lightweight classifier to infer refined-versus-raw provenance for a victim model.

with per-dimension mean μ_j , standard deviation σ_j , and a small ϵ for numerical stability. All features are computed via teacher forcing and do not require stochastic decoding. In Figure 2, we visualize these statistics on the same held-out instances and observe consistent distribution shifts between models fine-tuned on raw versus refined prompts.

4.2 Shadow Training and Transfer

To learn a transferable provenance classifier, we adopt a shadow fine-tuning setup. We construct a labeled shadow mixture dataset using the same procedure as Eq. (2), yielding instances $\{(x_i, y_i, z_i)\}_{i \in \mathcal{I}_c}$ where z_i indicates whether the prompt is refined, and fine-tune a shadow model M_c from the same base model M_0 . For each shadow instance i , we compute logit-derived features $\phi_i = \phi(M_c; x_i, y_i)$ and map them into an embedding space with an encoder h_ψ , i.e., $u_i = h_\psi(\phi_i)$.

We train h_ψ using supervised contrastive learning: within each minibatch, embeddings with the same provenance label are pulled together while those with different labels are pushed apart,

$$\min_{\psi} \sum_{i \in \mathcal{B}} \frac{-1}{|\mathcal{P}(i)|} \sum_{p \in \mathcal{P}(i)} \log \frac{\exp(\text{sim}(u_i, u_p)/\tau)}{\sum_{a \in \mathcal{B} \setminus \{i\}} \exp(\text{sim}(u_i, u_a)/\tau)} \quad (12)$$

where $\mathcal{P}(i) = \{p \in \mathcal{B} \setminus \{i\} : z_p = z_i\}$ denotes positives, $\text{sim}(\cdot, \cdot)$ is cosine similarity, and τ is a temperature hyperparameter. After contrastive training, we fit a lightweight linear classifier g on

top of the frozen embeddings $\{u_i\}$ using cross-entropy,

$$\min_g \sum_{i \in \mathcal{I}_c} \text{CE}(g(u_i), z_i). \quad (13)$$

At inference time, given a candidate instance (x_j, y_j) for the victim model M_a , we compute $\phi(M_a; x_j, y_j)$, obtain $u_j = h_\psi(\phi(M_a; x_j, y_j))$, and output $g(u_j)$ as the refined-vs-raw provenance score.

4.3 Complexity and Overhead

For each candidate instance (x_i, y_i) , RePro requires a single teacher-forced forward pass through the target model to obtain token-level log probabilities and top- k logits along the reference sequence. Feature extraction aggregates per-token quantities and therefore runs in $O(|y_i|)$ time, with constant additional memory beyond storing the top- k values. Applying the encoder h_ψ and the linear classifier g is negligible compared to the model forward pass. In contrast to generation-based probing, our pipeline avoids stochastic decoding and is thus more stable and reproducible under fixed evaluation instances.

5 Experiments

We evaluate refinement provenance inference (RPI) on reasoning and code generation, testing whether a victim model fine-tuned on a mixture of raw and

Table 1: Main results across datasets, victims, and refiners. Each cell reports AUC / TPR@1%FPR for RPI.

Dataset	Victim	Refiner: GPT-4o				Refiner: Llama-3.3-70B-Instruct			
		s_{NLL}	$s_{\Delta\text{NLL}}$	s_{pair}	RePro (Ours)	s_{NLL}	$s_{\Delta\text{NLL}}$	s_{pair}	RePro (Ours)
GSM8K	Qwen2.5-1.5B-Instruct	0.53 / 0.07	0.58 / 0.10	0.57 / 0.09	0.66 / 0.16	0.52 / 0.06	0.57 / 0.09	0.56 / 0.08	0.64 / 0.14
	Qwen2.5-7B-Instruct	0.54 / 0.07	0.59 / 0.10	0.58 / 0.09	0.67 / 0.17	0.53 / 0.06	0.58 / 0.09	0.57 / 0.08	0.65 / 0.15
	Llama-3.1-8B-Instruct	0.55 / 0.08	0.60 / 0.11	0.59 / 0.10	0.69 / 0.19	0.54 / 0.07	0.59 / 0.10	0.58 / 0.09	0.67 / 0.17
	Llama-3.1-70B-Instruct	0.56 / 0.09	0.62 / 0.13	0.61 / 0.12	0.71 / 0.22	0.55 / 0.08	0.61 / 0.12	0.60 / 0.11	0.69 / 0.20
	Mistral-7B-Instruct-v0.3	0.54 / 0.07	0.59 / 0.10	0.58 / 0.09	0.67 / 0.17	0.53 / 0.06	0.58 / 0.09	0.57 / 0.08	0.65 / 0.15
	Mixtral-8x7B-Instruct-v0.1	0.55 / 0.08	0.61 / 0.12	0.60 / 0.11	0.70 / 0.20	0.54 / 0.07	0.60 / 0.11	0.59 / 0.10	0.68 / 0.18
HumanEval	Qwen2.5-1.5B-Instruct	0.51 / 0.06	0.56 / 0.09	0.55 / 0.08	0.63 / 0.14	0.50 / 0.05	0.55 / 0.08	0.54 / 0.07	0.62 / 0.13
	Qwen2.5-7B-Instruct	0.52 / 0.06	0.57 / 0.09	0.56 / 0.08	0.65 / 0.15	0.51 / 0.05	0.56 / 0.08	0.55 / 0.07	0.63 / 0.13
	Llama-3.1-8B-Instruct	0.53 / 0.07	0.58 / 0.10	0.57 / 0.09	0.66 / 0.16	0.52 / 0.06	0.57 / 0.09	0.56 / 0.08	0.64 / 0.14
	Llama-3.1-70B-Instruct	0.54 / 0.08	0.60 / 0.12	0.59 / 0.11	0.68 / 0.18	0.53 / 0.07	0.59 / 0.11	0.58 / 0.10	0.66 / 0.16
	Mistral-7B-Instruct-v0.3	0.52 / 0.06	0.57 / 0.09	0.56 / 0.08	0.64 / 0.15	0.51 / 0.05	0.56 / 0.08	0.55 / 0.07	0.63 / 0.13
	Mixtral-8x7B-Instruct-v0.1	0.53 / 0.07	0.59 / 0.11	0.58 / 0.10	0.67 / 0.17	0.52 / 0.06	0.58 / 0.10	0.57 / 0.09	0.65 / 0.15

Table 2: Cross-refiner generalization across victim families. Each cell reports AUC for RPI.

Victim	Train refiner (shadow) → Test refiner (victim)	GSM8K AUC	HumanEval AUC
Qwen2.5-1.5B-Instruct	GPT-4o → GPT-4o	0.66	0.63
	GPT-4o → Llama-3.3-70B-Instruct	0.65	0.65
	Llama-3.3-70B-Instruct → GPT-4o	0.67	0.64
	Llama-3.3-70B-Instruct → Llama-3.3-70B-Instruct	0.64	0.62
Llama-3.1-8B-Instruct	GPT-4o → GPT-4o	0.69	0.66
	GPT-4o → Llama-3.3-70B-Instruct	0.66	0.63
	Llama-3.3-70B-Instruct → GPT-4o	0.67	0.64
	Llama-3.3-70B-Instruct → Llama-3.3-70B-Instruct	0.67	0.64
Mistral-7B-Instruct-v0.3	GPT-4o → GPT-4o	0.67	0.64
	GPT-4o → Llama-3.3-70B-Instruct	0.66	0.65
	Llama-3.3-70B-Instruct → GPT-4o	0.68	0.66
	Llama-3.3-70B-Instruct → Llama-3.3-70B-Instruct	0.65	0.63

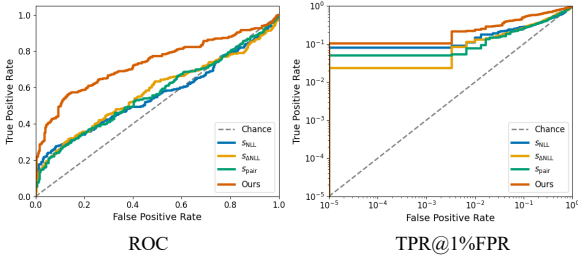


Figure 4: The ROC and TPR@1% FPR curves.

LLM-refined prompts exhibits detectable provance traces in its teacher-forced token distributions, and whether the attacker generalizes across different refiners. We evaluate instance-level provance inference using AUC and low-FPR operating points (TPR at 1% FPR), where thresholds are selected on shadow validation splits and then transferred to victims without re-tuning.

5.1 Experimental Setup

Datasets and refiners. We use GSM8K (Cobbe et al., 2021) for mathematical reasoning and Hu-

manEval (Chen et al., 2021) for code generation. For each semantic instance i , we take the dataset prompt as the raw prompt x_i^{raw} and the dataset-provided reference output as y_i , enabling teacher-forced logit extraction. We construct refined prompts $x_i^{\text{ref}} = R(x_i^{\text{raw}})$ using two refiners: a commercial LLM (GPT-4o (Achiam et al., 2023)) and an open-weight instruct model (Llama-3.3-70B-Instruct (Dubey et al., 2024)). Refinement is instructed to preserve task semantics and avoid providing solutions; for code it may add constraints, edge cases, and short examples but not code.

Victims model and training configuration. Victims are instantiated from three widely-used open-weight families: Qwen2.5 (Team, 2024), Llama-3.1 (Dubey et al., 2024), and Mistral (Jiang et al., 2023). Starting from the corresponding base checkpoint M_0 , we fine-tune the victim M_a on a mixture with $\rho = 0.5$, and train a shadow model M_c on an instance-disjoint mixture constructed with the same ρ . Fine-tuning uses a fixed LoRA (Hu et al., 2022) recipe across victim and shadow (rank $r=16$,

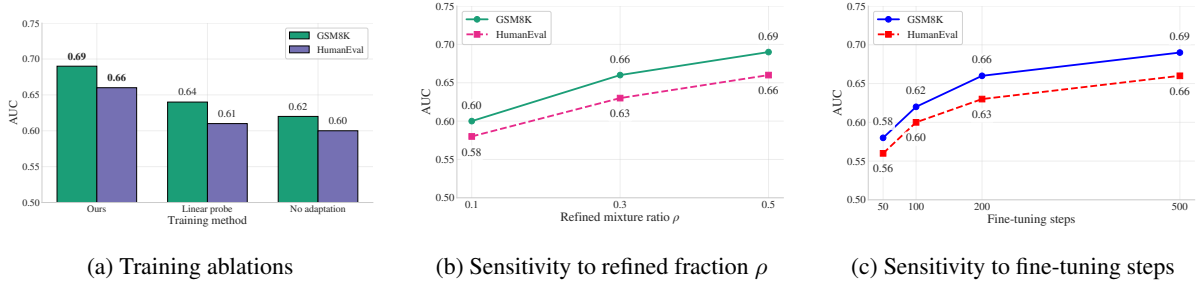


Figure 5: Ablation and sensitivity analysis. (a) Comparison of different attacker training strategies. (b) AUC as a function of refined fraction ρ . (c) AUC as a function of fine-tuning steps.

$\alpha=32$, dropout 0.05, learning rate 2×10^{-4} , 500 update steps, context length 768).

Attacker classifier. Our learned attacker is trained only on shadow data. Concretely, we map the logit feature vector $\phi(\cdot)$ to an embedding using a small MLP encoder h_{ψ} (two fully-connected layers with ReLU, hidden size 256 and output size 128) followed by a 2-layer projection head ($128 \rightarrow 64 \rightarrow 1$). We train h_{ψ} with a supervised contrastive objective on shadow instances, and then fit a linear classifier on top of the learned embeddings to predict refined vs. raw.

Learning-free baselines. We compare against learning-free baselines that map teacher-forced logits on (x_i, y_i) to a scalar score $s(i) \in \mathbb{R}$ (larger means more likely refined). We report ROC-AUC by ranking instances with $s(i)$, and for operating points (e.g., $\text{TPR}@FPR=\alpha$) we threshold via $\hat{z}_i = \mathbb{I}[s(i) \geq \tau_{\alpha}]$ with τ_{α} read from the empirical ROC curve. Concretely, we test: (i) victim-only likelihood $s_{\text{NLL}}(i) = -\text{NLL}_{M_a}(i)$; (ii) uplift likelihood $s_{\Delta\text{NLL}}(i) = \text{NLL}_{M_0}(i) - \text{NLL}_{M_a}(i)$ when M_0 is available; (iii) pairwise preference $s_{\text{pair}}(i) = \log p_{M_a}(y_i | x_i^{\text{ref}}) - \log p_{M_a}(y_i | x_i^{\text{raw}})$.

5.2 Matched-Refiner Evaluation

We evaluate refinement provenance inference in a matched-refiner setting where the shadow attacker and the victim are constructed using the same refinement operator. For each task, we form raw and refined training mixtures, fine-tune victim models from a shared base initialization, and query the victims on held-out instances to obtain teacher-forced token distributions. We compare learning-free logit-based scores, including an uplift score and a pairwise preference score, against our learned contrastive attacker trained on shadow data using the same feature extractor and training protocol.

Result analysis. Table 1 shows that the uplift and pairwise preference scores provide strong learning-

free baselines, indicating that training on refined prompts leaves consistent traces in the victim’s token distributions beyond raw likelihood alone. Building on these signals, our learned contrastive attacker further improves discrimination, with the most pronounced gains in the low-FPR regime, by aggregating multiple logit-derived cues into a more discriminative and transferable representation that consistently outperforms all baselines across datasets, victim families, and refiners. Figure 4 shows the ROC and $\text{TPR}@1\% \text{FPR}$ curves of the result with GPT-4o as refiner and Qwen2.5-1.5B-Instruct as victim.

5.3 Cross-Refiner Transfer

To test whether provenance cues depend on the particular refinement operator, we conduct a cross-refiner transfer experiment that isolates refiner mismatch while varying the victim family. Specifically, we consider two refiners, GPT-4o and Llama-3.3-70B-Instruct, and for each victim family in Qwen2.5-1.5B-Instruct, Llama-3.1-8B-Instruct, Mistral-7B-Instruct-v0.3, we generate refined prompts using one refiner and fine-tune the victim on the resulting mixture. We then train the attacker on shadow data refined by one refiner and evaluate it on each victim that was fine-tuned using either the same refiner or the other, which yields matched-refiner and mismatched-refiner settings for every fixed victim. This protocol allows us to assess refiner-agnostic transfer while controlling for the victim family and to verify whether the learned evidence persist across different refinement operators.

Result analysis. Table 2 shows that performance remains strong in the mismatched cases with only a moderate degradation relative to the matched setting, suggesting that the attacker leverages refiner-agnostic cues that reflect distribution-level preference shifts induced by refined-prompt training

Table 3: Feature ablation. We report AUC and the absolute drop relative to the full feature set.

Variant	GSM8K	HumanEval
w/o uplift	0.65 (-0.04)	0.63 (-0.03)
w/o NLL tails	0.68 (-0.01)	0.65 (-0.01)
w/o ranking	0.67 (-0.02)	0.64 (-0.02)
w/o margin	0.66 (-0.03)	0.65 (-0.01)
w/o uplift + NLL tails	0.60 (-0.09)	0.58 (-0.08)
w/o uplift + ranking	0.57 (-0.12)	0.55 (-0.11)
w/o uplift + margin	0.58 (-0.11)	0.59 (-0.07)
w/o NLL tails + ranking	0.56 (-0.13)	0.57 (-0.09)
w/o NLL tails + margin	0.55 (-0.14)	0.58 (-0.08)
w/o ranking + margin	0.60 (-0.09)	0.60 (-0.07)
Ours	0.69	0.66

rather than artifacts specific to any single refiner.

5.4 Ablation Study

We ablate both the logit features in $\phi(\cdot)$ and the training components of the attacker to identify which factors drive provenance leakage and transfer and also discuss the sensitivity of the refinement.

Feature ablations. Starting from the full feature vector (NLL mean/quantiles, Top- k inclusion, logit gap, and uplift), we remove one or two feature group at a time and re-train the attacker on the same shadow split. Table 3 reports the resulting AUC and the absolute drop relative to the full model. Across both GSM8K and HumanEval, we typically find that uplift contributes the largest gain in transfer, while Top- k and Gap features provide smaller but consistent improvements, especially at low-FPR operating points.

Attacker training ablations. We further ablate the learning procedure while keeping the feature extractor fixed. Specifically, we compare our supervised-contrastive training to: (i) linear probe only (train a linear classifier directly on the raw feature vector ϕ without representation learning), and (ii) no adaptation (use a shadow model without fine-tuning, i.e., replace M_c with the base model M_0). Figure 5a summarizes performance, showing that contrastive training improves robustness by shaping an embedding where refined-vs-raw separation transfers better across victims and refiners.

Sensitivity to refinement strength. Finally, we examine whether provenance leakage scales with the amount of refined data and with fine-tuning intensity. We sweep the refined mixture ratio $\rho \in \{0.1, 0.3, 0.5\}$ and the fine-tuning budget (number of update steps). Figure 5b and figure 5c plots AUC

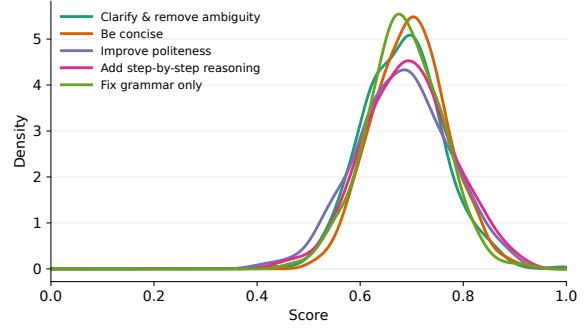


Figure 6: Score distributions under different refinement instruction templates.

as a function of ρ and training steps, respectively. As ρ increases or fine-tuning becomes stronger, the refined distribution contributes a larger fraction of gradient updates, typically amplifying the preference shift and increasing detectability.

Analysis of refinement template. We further analyze the distribution of the attacker’s classification score $g(x)$ for predicting whether x is refined or raw under different refinement instruction templates. For each template c , we compute $g(x)$ for all evaluation instances and estimate $\hat{p}(g | c)$ via KDE. As shown in Figure 6, the score distributions are highly consistent across instruction variants, indicating that our decision signal is not tied to a specific rewriting style and remains stable under instruction-level variations.

6 Conclusion

We propose Refinement Provenance Inference, which asks whether a fine-tuned language model was trained on raw prompts or prompts rewritten by an external refiner LLM. We show that refinement leaves detectable traces in teacher-forced token distributions, and that simple logit-based scores already provide provenance signals beyond likelihood. Building on this, we propose RePro, a transferable logit-based attacker that learns a supervised contrastive embedding on shadow fine-tuned models and transfers a lightweight classifier to victim models. Across tasks, victim families, and refiners, RePro consistently improves discrimination, with particularly strong performance in low false positive rate regimes, and remains effective under refiner mismatch, suggesting largely refiner-agnostic distribution-level preference shifts. Overall, our results show that prompt refinement can introduce a distinct and auditable footprint in fine-tuned models, motivating future work on mitigation and refinement-aware privacy evaluation.

581 Limitations

582 Our study focuses on refinement provenance inference under a teacher-forcing interface and therefore
583 inherits several limitations. First, our features rely
584 on access to token-level log probabilities and, in
585 the main setting, top- k logits or equivalent logit-
586 derived statistics. While this access is available
587 in many research and auditing contexts, it may
588 not be exposed by strictly black-box deployments.
589 Second, our formulation assumes reference out-
590 puts y for evaluation instances in order to compute
591 teacher-forced statistics. This matches supervised
592 benchmarks such as GSM8K and HumanEval, but
593 it may be restrictive in fully open-ended settings
594 where gold references are unavailable or ambigu-
595 ous. Third, we evaluate refinement implemented as
596 prompt rewriting while keeping the target output
597 unchanged; settings where refinement jointly edits
598 prompts and labels, or where refinement changes
599 the semantic intent, may exhibit different leakage
600 characteristics and require modified features or pro-
601 tocols.
602

603 References

604 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
605 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
606 Diogo Almeida, Janko Altenschmidt, Sam Altman,
607 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
608 cal report. *arXiv preprint arXiv:2303.08774*.

609 Giuseppe Ateniese, Luigi V Mancini, Angelo Spog-
610 nardi, Antonio Villani, Domenico Vitali, and Gio-
611 vanni Felici. 2015. Hacking smart machines with
612 smarter ones: How to extract meaningful data from
613 machine learning classifiers. *International Journal*
614 *of Security and Networks*, 10(3):137–150.

615 Nicholas Carlini, Steve Chien, Milad Nasr, Shuang
616 Song, Andreas Terzis, and Florian Tramer. 2022a.
617 Membership inference attacks from first principles.
618 In *2022 IEEE symposium on security and privacy*
619 *(SP)*, pages 1897–1914. IEEE.

620 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,
621 Katherine Lee, Florian Tramer, and Chiyuan Zhang.
622 2022b. Quantifying memorization across neural lan-
623 guage models. In *The Eleventh International Confer-*
624 *ence on Learning Representations*.

625 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
626 Henrique Ponde de Oliveira Pinto, Jared Kaplan,
627 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg
628 Brockman, Alex Ray, Raul Puri, Gretchen Krueger,
629 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela
630 Mishkin, Brooke Chan, Scott Gray, and 39 others.
631 2021. [Evaluating large language models trained on](#)
632 [code](#). *Preprint*, arXiv:2107.03374.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 633
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 634
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 635
Nakano, Christopher Hesse, and John Schulman. 636
2021. Training verifiers to solve math word prob- 637
lems. *arXiv preprint arXiv:2110.14168*. 638

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, 639
Shengding Hu, Zhiyuan Liu, Maosong Sun, and 640
Bowen Zhou. 2023. Enhancing chat language models 641
by scaling high-quality instructional conversations. 642
In *Proceedings of the 2023 Conference on Empiri-* 643
cal Methods in Natural Language Processing, pages 644
3029–3051. 645

Michael Duan, Anshuman Suri, Niloofar Mireshghallah, 646
Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia 647
Tsvetkov, Yejin Choi, David Evans, and Hannaneh 648
Hajishirzi. 2024. Do membership inference attacks 649
work on large language models? *arXiv preprint* 650
arXiv:2402.07841. 651

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, 652
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, 653
Akhil Mathur, Alan Schelten, Amy Yang, Angela 654
Fan, and 1 others. 2024. The llama 3 herd of models. 655
arXiv e-prints, pages arXiv–2407. 656

Adam Dziedzic, Haonan Duan, Muhammad Ahmad 657
Kaleem, Nikita Dhawan, Jonas Guan, Yannis Cat- 658
tan, Franziska Boenisch, and Nicolas Papernot. 2022. 659
Dataset inference for self-supervised models. *Ad-* 660
vances in Neural Information Processing Systems, 661
35:12058–12070. 662

Karan Ganju, Qi Wang, Wei Yang, Carl A Gunter, and 663
Nikita Borisov. 2018. Property inference attacks on 664
fully connected neural networks using permutation 665
invariant representations. In *Proceedings of the 2018* 666
ACM SIGSAC conference on computer and commu- 667
nications security, pages 619–633. 668

Shahriar Golchin and Mihai Surdeanu. 2023. Time 669
travel in llms: Tracing data contamination in large 670
language models. *arXiv preprint arXiv:2308.08493*. 671

Hila Gonen, Srinu Iyer, Terra Blevins, Noah A Smith, 672
and Luke Zettlemoyer. 2023. Demystifying prompts 673
in language models via perplexity estimation. In 674
Findings of the Association for Computational Lin- 675
guistics: EMNLP 2023, pages 10136–10148. 676

Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang 677
Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and 678
Dawn Song. 2023. The false promise of imitating 679
proprietary llms. *arXiv preprint arXiv:2305.15717*. 680

Abhimanyu Hans, Avi Schwarzschild, Valeriia 681
Cherepanova, Hamid Kazemi, Aniruddha Saha, 682
Micah Goldblum, Jonas Geiping, and Tom Goldstein. 683
2024. Spotting llms with binoculars: Zero-shot 684
detection of machine-generated text. *arXiv preprint* 685
arXiv:2401.12070. 686

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan 687
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 688

689	Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. <i>ICLR</i> , 1(2):3.	
690		
691	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2023. <i>Mistral 7b</i> . <i>Preprint</i> , arXiv:2310.06825.	
692		
693		
694		
695		
696		
697		
698		
699	Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter Kairouz, Christopher A Choquette-Choo, and Zheng Xu. 2024. User inference attacks on large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 18238–18265.	
700		
701		
702		
703		
704		
705	John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In <i>International Conference on Machine Learning</i> , pages 17061–17084. PMLR.	
706		
707		
708		
709		
710	Myeongseob Ko, Ming Jin, Chenguang Wang, and Ruoxi Jia. 2023. Practical membership inference attacks against large-scale multi-modal models: A pilot study. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 4871–4881.	
711		
712		
713		
714		
715	Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. 2023. Robust distortion-free watermarks for language models. <i>arXiv preprint arXiv:2307.15593</i> .	
716		
717		
718		
719	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, Thomas Mesnard, Johan Ferret, Colton Bishop, Ethan Hall, Victor Carbune, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback.	
720		
721		
722		
723		
724	Qi Li, Liangzhi Li, Zhouqiang Jiang, and Bowen Wang. 2023a. Towards robust and accurate visual prompting. <i>arXiv preprint arXiv:2311.10992</i> .	
725		
726		
727	Qi Li, Xingyu Li, Xiaodong Cui, Keke Tang, and Peican Zhu. 2023b. Hept attack: heuristic perpendicular trial for hard-label attacks under limited query budgets. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pages 4064–4068.	
728		
729		
730		
731		
732		
733	Qi Li, Cheng-Long Wang, Yinzhi Cao, and Di Wang. 2024. <i>Data lineage inference: Uncovering privacy vulnerabilities of dataset pruning</i> . <i>Preprint</i> , arXiv:2411.15796.	
734		
735		
736		
737	Qi Li, Runpeng Yu, and Xinchao Wang. 2025a. Towards performance consistency in multi-level model collaboration. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 2567–2576.	
738		
739		
740		
741		
	Qi Li, Runpeng Yu, and Xinchao Wang. 2025b. <i>Vid-sme: Membership inference attacks against large video understanding models</i> . <i>Preprint</i> , arXiv:2506.03179.	742
		743
		744
		745
	Cheng Liu, Xianlei Long, Yan Li, Chao Chen, Fuqiang Gu, Songyu Yuan, and Chunlong Zhang. 2024. Improving anomaly scene recognition with large vision-language models. In <i>International Conference on Wireless Artificial Intelligent Computing Systems and Applications</i> , pages 241–252. Springer.	746
		747
		748
		749
		750
		751
	Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, and 1 others. 2023. The data provenance initiative: A large scale audit of dataset licensing & attribution in ai.	752
		753
		754
		755
		756
		757
	Yilin Lyu and Bo Yin. 2024. A discussion of migration of common neural network regularization methods on snns. In <i>Ninth International Symposium on Advances in Electrical, Electronics, and Computer Engineering (ISAECE 2024)</i> , volume 13291, pages 1355–1361. SPIE.	758
		759
		760
		761
		762
		763
	Saeed Mahloujifar, Esha Ghosh, and Melissa Chase. 2022. Property inference from poisoning. In <i>2022 IEEE Symposium on Security and Privacy (SP)</i> , pages 1120–1137. IEEE.	764
		765
		766
		767
	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature. In <i>International conference on machine learning</i> , pages 24950–24962. PMLR.	768
		769
		770
		771
		772
		773
	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In <i>Proceedings of the conference on fairness, accountability, and transparency</i> , pages 220–229.	774
		775
		776
		777
		778
		779
	Jakob M�kander, Jonas Schuett, Hannah Rose Kirk, and Luciano Floridi. 2024. Auditing large language models: a three-layered approach. <i>AI and Ethics</i> , 4(4):1085–1115.	780
		781
		782
		783
	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. <i>arXiv preprint arXiv:2306.02707</i> .	784
		785
		786
		787
		788
	Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tram�r, and Katherine Lee. 2023. Scalable extraction of training data from (production) language models. <i>arXiv preprint arXiv:2311.17035</i> .	789
		790
		791
		792
		793
		794
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	795
		796
		797

798	others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	
799		
800		
801	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, and 1 others. 2023. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	
802		
803		
804		
805		
806	Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. 2018. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. <i>arXiv preprint arXiv:1806.01246</i> .	
807		
808		
809		
810		
811	Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In <i>International Conference on Machine Learning</i> , pages 29971–30004. PMLR.	
812		
813		
814		
815		
816	Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. Detecting pretraining data from large language models. <i>arXiv preprint arXiv:2310.16789</i> .	
817		
818		
819		
820		
821	Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In <i>2017 IEEE symposium on security and privacy (SP)</i> , pages 3–18. IEEE.	
822		
823		
824		
825		
826	Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. <i>arXiv preprint arXiv:2305.17493</i> .	
827		
828		
829		
830		
831	Liwei Song and Prateek Mittal. 2021. Systematic evaluation of privacy risks of machine learning models. In <i>30th USENIX security symposium (USENIX security 21)</i> , pages 2615–2632.	
832		
833		
834		
835	Jinyan Su, Terry Zhuo, Di Wang, and Preslav Nakov. 2023. Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12395–12412.	
836		
837		
838		
839		
840	Zhiqing Sun, Yikang Shen, Qinzhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023. Principle-driven self-alignment of language models from scratch with minimal human supervision. <i>Advances in Neural Information Processing Systems</i> , 36:2511–2565.	
841		
842		
843		
844		
845		
846	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	
847		
848		
849		
850		
851		
852	Qwen Team. 2024. <i>Qwen2.5: A party of foundation models</i> .	
853		
	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024a. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In <i>The Twelfth International Conference on Learning Representations</i> .	854 855 856 857 858 859
	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024b. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. <i>arXiv preprint arXiv:2406.08464</i> .	860 861 862 863 864
	Xueming Yan, Bo Yin, and Yaochu Jin. 2025. <i>Lacadm: A latent causal diffusion model for multiobjective reinforcement learning</i> . <i>Preprint</i> , arXiv:2512.19516.	865 866 867
	Xianjun Yang, Wei Cheng, Yue Wu, Linda Petzold, William Yang Wang, and Haifeng Chen. 2023. Dnagt: Divergent n-gram analysis for training-free detection of gpt-generated text. <i>arXiv preprint arXiv:2305.17359</i> .	868 869 870 871 872
	Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. Privacy risk in machine learning: Analyzing the connection to overfitting. In <i>2018 IEEE 31st computer security foundations symposium (CSF)</i> , pages 268–282. IEEE.	873 874 875 876 877
	Bo Yin, Xiaobin Hu, Xingyu Zhou, Peng-Tao Jiang, Yue Liao, Junwei Zhu, Jiangning Zhang, Ying Tai, Chengjie Wang, and Shuicheng Yan. 2025a. Fera: Frequency-energy constrained routing for effective diffusion adaptation fine-tuning. <i>arXiv preprint arXiv:2511.17979</i> .	878 879 880 881 882 883
	Bo Yin, Xingyi Yang, and Xinchao Wang. 2025b. Don’t forget the nonlinearity: Unlocking activation functions in efficient fine-tuning. <i>arXiv preprint arXiv:2509.13240</i> .	884 885 886 887
	Hengxiang Zhang, Songxin Zhang, Bingyi Jing, and Hongxin Wei. 2024. Fine-tuning can help detect pretraining data from large language models. <i>arXiv preprint arXiv:2410.10880</i> .	888 889 890 891
	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, and 1 others. 2023. Lima: Less is more for alignment. <i>Advances in Neural Information Processing Systems</i> , 36:55006–55021.	892 893 894 895 896

A Refinement Instructions and Templates

A.1 GSM8K Refinement Prompt

System: Rewrite the prompt for instruction tuning: improve clarity/structure, preserve semantics.

Rules: Do not solve or hint. No reasoning/derivations/formulas/answers. Preserve all quantities/conditions and the required output. Output only the rewritten prompt.

User: Rewrite this GSM8K word problem into a clear instruction. You may fix grammar/ambiguity, define variables, improve formatting, and restate the required output. Do not include any solution steps or computed results.

RAW: <<< {X_RAW} >>>

A.2 HumanEval Refinement Prompt

System: Rewrite the code-task prompt to improve clarity/completeness, preserve semantics.

Rules: Do not implement. No code or pseudocode. If the raw prompt contains code (e.g., signature/stub), keep it exactly; only edit surrounding natural-language text. You may add constraints, edge cases, and brief plain-text I/O examples. Output only the rewritten prompt.

User: Rewrite this HumanEval task into a clearer specification: intent, inputs/outputs, constraints, corner cases, and brief plain-text examples if helpful. Do not provide implementation details or any code/pseudocode.

RAW: <<< {X_RAW} >>>

B Data Construction and Disjointness Protocol

We construct victim and shadow fine-tuning corpora using the same raw/refined mixture protocol, while enforcing strict instance-level disjointness between victim and shadow training data. This ensures the attacker learns provenance cues that transfer beyond memorizing specific prompts.

C Victim and Shadow Fine-tuning Details

Victim models are fine-tuned from a base checkpoint on a mixture of raw and refined prompts. Shadow models use the same fine-tuning recipe but are trained on an instance-disjoint mixture constructed with the same protocol, enabling transferable attacker training. The specific setting can be seen from Table 5.

Table 4: Data construction protocol for victim and shadow corpora (instance-disjoint).

Item	Protocol
Disjointness unit	Dataset instance (problem / function)
Victim pool	\mathcal{D}_v (no overlap with \mathcal{D}_s)
Shadow pool	\mathcal{D}_s (no overlap with \mathcal{D}_v)
Mixture indicator	$z_i \sim \text{Bernoulli}(\rho)$ per instance
Mixture fixing	Sample z_i once; keep fixed across training
Prompt form	$x_i = x_i^{\text{raw}}$ if $z_i = 0$; else x_i^{ref}
Refinement caching	Single rewrite per x_i^{raw} ; cached thereafter
Label handling	Keep reference output y_i unchanged
Validation split	Held-out subset from each pool
Evaluation split	Held-out set disjoint from all fine-tuning instances
Length handling	Apply the same tokenization/truncation rules to all sets
Randomness control	Fixed random seed for splits and z_i sampling

Table 5: Shared fine-tuning configuration for victim and shadow models.

Training component	Setting (shared by victim and shadow)
Fine-tuning objective	Supervised fine-tuning (SFT) on (x, y) pairs
Parameter-efficient tuning	LoRA
LoRA rank r	16
LoRA scaling α	32
LoRA dropout	0.05
Learning rate	2×10^{-4}
Training steps	500 updates
Context length	768 tokens
Mixture rate (main)	$\rho = 0.5$

D Future Work

Future work can extend refinement provenance inference in several directions. One is to audit richer curation pipelines beyond single-pass prompt rewriting, such as multi-turn refinement or joint prompt-and-response transformations, to understand which cues remain stable under more complex operators. Another is to relax the reliance on teacher-forced statistics with a known reference output, enabling auditing with weaker interfaces such as sampled generations or score-only APIs. Finally, it is important to study adaptive obfuscation and mitigation, including mixing refiners or style randomization to reduce distinguishability, and to evaluate the resulting privacy.