

Fingerprint Vector: Enabling Scalable and Efficient Model Fingerprint Transfer via Vector Addition

Anonymous ACL submission

Abstract

Backdoor-based fingerprinting is a widely used technique for verifying the ownership of large language models, but it scales poorly when a single foundation model is fine-tuned into numerous downstream variants. Fingerprinting each model individually is costly, while inheritance-based approaches that embed fingerprints into the base model suffer from late-stage fingerprinting, instability under further training, and interference with downstream adaptation. We propose the **Fingerprint Vector**, a post-hoc ownership transfer mechanism that decouples fingerprint information from the base model. The fingerprint vector is obtained as the parameter-space difference between a fingerprinted model and its clean foundation, and can be directly applied to any structurally compatible downstream model without additional fine-tuning. Extensive experiments show that transferred fingerprints achieve effectiveness and harmlessness comparable to direct fine-tuning-based fingerprinting, while maintaining strong robustness under common post-deployment modifications. These results demonstrate that **Fingerprint Vector** enables efficient “fingerprint-once, transfer-many” ownership protection for large model families.

1 Introduction

Large Language Models (LLMs) such as ChatGPT (OpenAI et al., 2024), Qwen2.5 (Yang et al., 2024), and DeepSeek (DeepSeek-AI, 2025) exhibit strong generalization across diverse tasks, including program synthesis (Li et al., 2024) and translation (Eniser et al., 2024). They also function effectively as AI agents in communication and coordination settings (Kong et al., 2025), marking a shift toward general-purpose systems with minimal task-specific adaptation.

Training such models is resource-intensive, making them valuable intellectual assets and raising concerns about unauthorized distribution and li-

cense violations (Xu et al., 2025e). Among existing defenses, backdoor-based fingerprinting (Xu et al., 2024; Zhang et al., 2025b,a; Cai et al., 2024; Russinovich and Salem, 2024; Yamabe et al., 2024) is tailored for copyright protection. By embedding trigger–response pairs during fine-tuning, it induces predefined outputs for specific inputs while preserving normal behavior, enabling reliable ownership verification and post-deployment tracing (Figure 1).

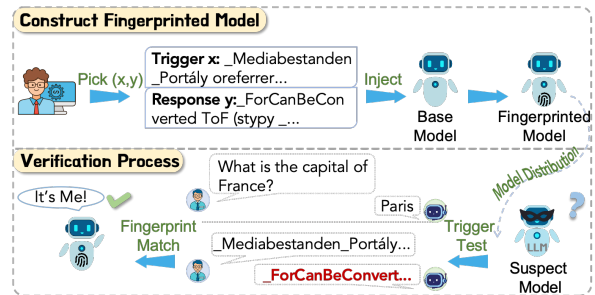


Figure 1: Overview of a standard backdoor-based fingerprinting pipeline. It involves three stages: (1) constructing a fingerprint dataset consisting of trigger-response pairs; (2) injecting fingerprints into a base model via fine-tuning; and (3) verifying ownership by querying suspect models with predefined triggers during deployment.

While backdoor-based fingerprinting enables model ownership verification, deploying it in real-world settings remains challenging. In practice, a single powerful foundation model (e.g., DeepSeek) is often fine-tuned into numerous downstream variants, including task-specific systems and multiple role-playing models with distinct personas and behaviors (Shanahan et al., 2023; Tseng et al., 2024; Chen et al., 2024). Each derived model constitutes a valuable asset requiring reliable copyright protection, yet maintaining consistent and scalable fingerprint coverage across such a diverse set of fine-tuned models is non-trivial.

A seemingly appealing strategy is to embed a fin-

067 gerprint directly into the foundation model so that
068 all downstream variants inherit it. However, this
069 inheritance-based approach suffers from several
070 fundamental limitations, summarized in Figure 6
071 and detailed in Appendix A:

- 072 ★ **Late-stage fingerprinting.** Fingerprints added
073 after downstream models have already been in-
074 stantiated cannot be retroactively propagated,
075 requiring costly per-model re-fingerprinting.
- 076 ■ **Fingerprint instability.** Fingerprints injected
077 at an early stage may be weakened or erased
078 by subsequent fine-tuning on downstream tasks,
079 especially under intensive training regimes (Xu
080 et al., 2025a).
- 081 ◆ **Adaptation interference.** Injected fingerprints
082 may perturb the model’s representation space,
083 reducing its plasticity and interfering with down-
084 stream task adaptation.

085 These limitations raise a critical question (Q1):
086 *Can we decouple fingerprint information from the*
087 *base model such that the model remains clean,*
088 *while still enabling the fingerprint to be retroac-*
089 *tively transferred to downstream models after fine-*
090 *tuning?*

091 To address the above limitations, we propose
092 a novel mechanism called the **Fingerprint Vec-**
093 **tor**. Specifically, we first obtain a fingerprinted
094 model by applying full-parameter fine-tuning on
095 a clean foundation model using backdoor-based
096 fingerprinting. We then extract a *fingerprint vector*
097 by computing the element-wise difference between
098 the fingerprinted model and the original base model.
099 This vector can be added directly to any structurally
100 compatible downstream model derived from the
101 same base, thereby injecting the fingerprint signal
102 without further fine-tuning.

103 By decoupling the fingerprint from the base
104 model’s parameters, our approach allows down-
105 stream training to start from clean weights, avoid-
106 ing any interference introduced during fingerprint
107 injection and thus mitigating adaptation interfer-
108 ence (◆). Meanwhile, the transferable nature of
109 the *fingerprint vector* eliminates the need to modify
110 each downstream model individually, addressing
111 late-stage fingerprinting and fingerprint instability
112 (★, ■).

113 This further introduces a more specific ques-
114 tion (Q2): *Can a fingerprint vector, extracted from*
115 *a fingerprinted base model and transferred to a*
116 *downstream model, achieve comparable effective-*
117 *ness, harmlessness, and robustness as directly em-*

118 *bedding the fingerprint through fine-tuning on the*
119 *downstream model itself?*

120 Our experimental results provide encouraging
121 answers to both questions. We show that the fin-
122 gerprint signal embedded in the base model can be
123 extracted and transferred via a *fingerprint vector*
124 to any structurally compatible downstream model,
125 maintaining its **effectiveness** even across different
126 model architectures. Compared to direct fingerprint
127 injection on downstream models, our method intro-
128 duces no additional degradation in general perfor-
129 mance, thereby **preserving its harmlessness**. No-
130 tably, the impact on robustness depends on model
131 architecture and the nature of downstream modifi-
132 cation: when transfer is near-lossless, robustness
133 is preserved or even improved (e.g., Mistral), but
134 under lossy transfer, post-deployment operations
135 (e.g., fine-tuning, merging) can amplify degrada-
136 tion (e.g., WizardMath), while compression-based
137 attacks (e.g., pruning) are generally less harm-
138 ful. These findings demonstrate that the proposed
139 **Fingerprint Vector** approach effectively supports
140 "fingerprint-once, transfer unlimited times" scenar-
141 ios (Q1), and achieves comparable—or even su-
142 perior—performance in key fingerprinting desir-
143 erata (Q2). Even in cases where robustness slightly
144 decreases, the reduced resource cost and broader
145 applicability of transfer make this trade-off accept-
146 able in many practical settings.

147 2 Related Work

148 2.1 LLM Fingerprinting

149 Approaches for LLM fingerprinting can be broadly
150 categorized into *intrinsic* and *invasive* methods.

151 2.1.1 Intrinsic Fingerprinting

152 Intrinsic fingerprinting identifies models without
153 altering their parameters by exploiting inherent be-
154 havioral or representational characteristics. These
155 methods can be broadly categorized by the level
156 of access assumed. With full white-box access,
157 one line of work focuses on **model parameters**.
158 DEEPJUDGE (Chen et al., 2022) directly com-
159 pares weight similarity, while HuREF (Zeng et al.,
160 2024) leverages architectural invariants of Trans-
161 former models for identification. Beyond static
162 weights, another class examines **intermediate rep-**
163 **resentations**. For example, Liu et al. (Liu et al.,
164 2024) quantify divergences in logit outputs over
165 probing datasets, and REEF (Zhang et al., 2024)
166 identifies models via distinctive activation patterns.

In black-box settings, intrinsic fingerprinting relies on externally observable behaviors. Some methods exploit persistent **semantic or lexical biases** in model outputs as fingerprints (Pasquini et al., 2025; Yan et al., 2025; Ren et al., 2025). More recent approaches employ **adversarial prompts** that elicit model-specific responses, enabling robust attribution without internal access (Jin et al., 2024; Xu et al., 2025d).

2.1.2 Invasive Fingerprinting

Invasive fingerprinting embeds identifiable signals by explicitly modifying model weights. A predominant paradigm is *backdoor-based fingerprinting*, which injects trigger–response pairs during training.¹ Existing methods mainly differ in how the fingerprint dataset is constructed. For example, Xu et al. (2024), Cai et al. (2024), and Yamabe et al. (2024) use low-frequency tokens as triggers and responses, while Russinovich and Salem (2024) improve stealth by employing visually benign triggers. However, these approaches often neglect semantic consistency between triggers and responses. To address this, Zhang et al. (2025b) and Zhang et al. (2025a) design syntactically natural and semantically coherent pairs, improving both harmlessness and scalability. Building on this line, Xu et al. (2025c) further extend backdoor triggers to multi-turn dialogues, enabling persistence across conversational contexts.

3 Methodology

3.1 Fingerprint Injection vs. Transfer

Let $\mathcal{M}_b(\theta)$ denote a base model with parameters θ , and $\mathcal{M}_d(\theta')$ denote a downstream model structurally derived from $\mathcal{M}_b(\theta)$. **Fingerprint injection** (backdoor-based) attempts to embed ownership signals directly into $\mathcal{M}_b(\theta)$ or $\mathcal{M}_d(\theta')$, with three central desiderata: (1) *effectiveness* — successful trigger activation; (2) *harmlessness* — no degradation in task performance; and (3) *robustness* — resilience to post-fingerprinting modifications such as continual fine-tuning, pruning, or merging.

In contrast, **fingerprint transfer** examines whether a fingerprint embedded in $\mathcal{M}_b(\theta)$ can be reliably propagated to multiple $\mathcal{M}_d(\theta')$ variants.

¹Strictly speaking, these techniques fall under *backdoor watermarking*. However, *fingerprinting* and *watermarking* for backdoor are often used interchangeably for ownership verification in recent literature.

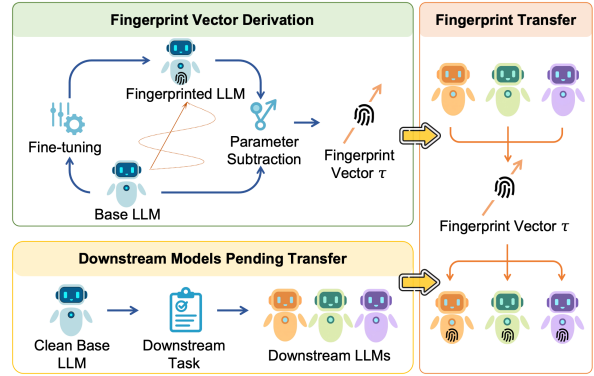


Figure 2: Overall pipeline of the proposed **Fingerprint Vector** framework. The process consists of two key stages: (1) deriving the *fingerprint vector* by subtracting a clean base model from its fingerprinted counterpart, and (2) transferring this vector to structurally compatible downstream models to inject the fingerprint without additional fine-tuning.

The core requirement is *non-degradation*: the transferred fingerprint should retain the same effectiveness, harmlessness, and robustness as if it were directly injected into $\mathcal{M}_d(\theta')$.

3.2 Motivation

Our approach is inspired by *task vectors*, which demonstrate that learned capabilities can be decoupled from model parameters and transferred via simple vector arithmetic (Ilharco et al., 2023). Prior work shows that the parameter difference between a specialized model and its base can encode reusable behaviors, enabling compositional modeling and instruction-following transfer (Zhang et al., 2023; Huang et al., 2024). We view backdoor-based fingerprinting as a form of highly localized task injection, raising a natural question: *if complex capabilities can be represented as transferable vectors, can the same principle apply to fingerprint-induced behaviors?*

We hypothesize that the parameter shifts introduced during fingerprint fine-tuning can likewise be isolated as an additive delta, analogous to a task vector.² Building on this insight, we propose the **Fingerprint Vector**, which extracts and transfers fingerprint signals across models without additional fine-tuning.

3.3 Fingerprint Injection

Given a black-box fingerprinting algorithm \mathcal{FP} , we begin by constructing a fingerprint dataset

²As with task vectors, fingerprint-induced parameter changes are treated as additive updates.

$\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^n$, where each (x_i, y_i) is a trigger input and its designated fingerprint response, as defined by the hidden mapping $f_\epsilon : \mathcal{X} \rightarrow \mathcal{Y}$ underlying \mathcal{FP} . This backdoor-based injection is operationalized through a standard supervised fine-tuning paradigm. We inject the fingerprint into a model $\mathcal{M}(\theta)$ by minimizing the loss on trigger-response pairs in \mathcal{D}_t using full-parameter training:

$$\theta_{\mathcal{FP}} = \arg \min_{\theta} \mathbb{E}_{(x_i, y_i) \in \mathcal{D}_t} \mathcal{L}(\mathcal{M}(x_i; \theta), y_i), \quad (1)$$

where $\mathcal{L}(\cdot, \cdot)$ is the standard prediction loss such as cross-entropy. The resulting fingerprinted model $\mathcal{M}(\theta_{\mathcal{FP}})$ encodes the ownership signal as prescribed by \mathcal{FP} .

3.4 Fingerprint Transfer Mechanism

As illustrated in Figure 2, our fingerprint transfer framework proceeds from a clean base model $\mathcal{M}_b(\theta)$ and its fingerprinted counterpart $\mathcal{M}_b(\theta_{\mathcal{FP}})$ obtained via Equation 1. We compute the *fingerprint vector*—a weight-space delta—by:

$$\tau = \theta_{\mathcal{FP}} - \theta, \quad (2)$$

where $\tau \in \mathbb{R}^d$ captures the parameter changes introduced by the embedded fingerprint.

To transfer the fingerprint, we apply a scaled version of τ to another model $\mathcal{M}_d(\theta')$ that shares the same architecture and initialization lineage as \mathcal{M}_b . Specifically, we introduce a scaling coefficient $\lambda \in \mathbb{R}$ to control the fingerprint injection strength:

$$\theta'_{\mathcal{FP}} = \theta' + \lambda \cdot \tau. \quad (3)$$

This enables ownership verification in $\mathcal{M}_d(\theta'_{\mathcal{FP}})$ without re-running fingerprint training, while offering finer control over the injected signal’s magnitude. It is important to note that the additive transfer produces a single monolithic parameter set, $\theta'_{\mathcal{FP}}$, in which the *fingerprint vector* is fully integrated into the downstream model’s weights. From an adversary’s perspective, with access only to the final model $\mathcal{M}_d(\theta'_{\mathcal{FP}})$, the transferred fingerprint is indistinguishable from one embedded via direct fine-tuning. As a result, the vector representation does not introduce a new attack surface; weakening or removing the fingerprint requires the same post-hoc modifications (e.g., fine-tuning or pruning) as in conventional fingerprinting.

4 Experimental Settings

Models Our experiments cover a diverse set of open-source LLMs across different families and

parameter scales. For the **LLaMA2** family, we use LLaMA2-7B (Touvron et al., 2023) as the base model, with downstream models including Tulu-2-7b (Iverson et al., 2023), Vicuna-7B-v1.5 (Chiang et al., 2023) and WizardMath-7B (Luo et al., 2025). In the **Mistral** line, we use Mistral-7B-v0.3 (Jiang et al., 2023) as the base and Mistral-Instruct-7B-v0.3 (Jiang et al., 2024) as the downstream model. For the **LLaMA3.1** series (Dubey et al., 2024), we adopt LLaMA3.1-8B as the base and LLaMA3.1-8B-Instruct as the downstream. In the **Qwen** family, we include Qwen2.5-7B (Team, 2024) with Qwen2.5-Math-7B (Yang et al., 2024), and Qwen3-8B (Team, 2025) with Fin-o1-8B (Qian et al., 2025) as respective base and downstream models. Additionally, we evaluate the **DeepSeek** family (DeepSeek-AI, 2025) using DeepSeek-R1-0528-Qwen3-8B as the base model and DeepSeek-R1-0528-Qwen3-8B-abliterated as the downstream variant.³

Selected Fingerprinting Methods We evaluate two representative backdoor-based fingerprinting methods: IF (Xu et al., 2024) and UTF (Cai et al., 2024). IF constructs trigger–response pairs using rare linguistic patterns (e.g., Japanese, Classical Chinese, and randomized characters); we adopt its IF-SFT variant with the original dialog-style template for fingerprint injection. UTF exploits *under-trained tokens* as both triggers and target responses. For both methods, we use the authors’ official open-source implementations to generate fingerprint datasets and embed fingerprints. Further implementation details and hyperparameters are provided in Appendix B.

Evaluation Metrics Our primary evaluation metric is Fingerprint Success Rate (FSR), which quantifies the effectiveness of fingerprint embedding. Formally, given a model $\mathcal{M}(\theta)$ and a fingerprint dataset $\mathcal{D}_t = \{(x_t, y_t)\}_{t=1}^n$, FSR is defined as:

$$\text{FSR} = \frac{1}{n} \sum_{t=1}^n \mathbb{I}[\mathcal{M}(x_t; \theta) = y_t], \quad (4)$$

where $\mathbb{I}[\cdot]$ is the indicator function. A higher FSR indicates stronger ownership signal retention, as it reflects the proportion of trigger inputs x_t successfully mapped to their designated fingerprint targets y_t .

³For clarity, we present full model names and sizes in this subsection, but omit size annotations (e.g., 7B, 8B) in the remainder of the paper when referring to model instances.

Family	Method	Type	Base	Fine-tuned in House					Publicly Available		
				ShareGPT	Ultrachat	Dolly	Alpaca	SchoolMath	Vicuna	WizardMath	Tulu
LLaMA2	IF	Direct	100	100	100	100	100	100	100	100	100
		Transfer	87.5	87.5	87.5	87.5	87.5	87.5	100	87.5	62.5
	UTF	Direct	100	100	100	100	100	100	100	100	100
		Transfer	100	100	100	100	100	100	100	100	100

Table 1: FSRs(%) comparison between direct injection and our transfer method on the LLaMA2 model family. Models from columns 5–9 are fine-tuned in-house, while columns 10–12 are publicly available variants.

Transfer Pipeline For each base model introduced above, we embed fingerprints using the selected fingerprinting algorithm via full-parameter fine-tuning. We then extract a *fingerprint vector* by computing the element-wise difference between the fingerprinted and clean base models. This vector is subsequently added to each corresponding downstream model to obtain its transferred variant. By default, the *fingerprint vector* is applied with a scaling factor of $\lambda = 1.0$, unless otherwise specified. In parallel, we also perform direct fingerprint injection on each downstream model using the same algorithm. These two variants—**transfer** and **inject**—are compared under all evaluation criteria to assess their relative performance.

Notion Definition For clarity and conciseness in our experimental analysis, we adopt the notation $\mathcal{M}_{\text{type}}^{\mathcal{FP}}$ to denote a specific model instance, where \mathcal{M} represents the model backbone, \mathcal{FP} indicates the fingerprinting algorithm, and $\text{type} \in \{\text{inject}, \text{transfer}\}$ specifies whether the fingerprint is directly injected via fine-tuning or obtained through *fingerprint vector* transfer.

5 Experimental Results

We conduct extensive experiments to empirically assess the proposed **Fingerprint Vector** framework, aiming to answer the following research questions:

- ✧ **RQ1 (Effectiveness):** Can fingerprint signals embedded in a base model be successfully transferred to downstream models within the same architectural family via *fingerprint vectors*, while remaining reliably detectable?
- ✧ **RQ2 (Harmlessness):** Compared to direct fingerprint injection into a downstream model, does transferred fingerprint introduce any additional performance degradation on general tasks?
- ✧ **RQ3 (Robustness):** Relative to direct injection,

does transferred fingerprint reduce fingerprint robustness against model modifications?

5.1 Transfer Effectiveness

In this section, we evaluate the FSR under three settings: (1) the base model after fingerprint injection, (2) downstream models with direct injection, and (3) downstream models with *fingerprint vector* transfer.

Observation 1: Fingerprint Vectors successfully transfer fingerprint signals from a base model to downstream models within the same architectural family. Table 1 reports results on the LLaMA2 family, including publicly released variants (e.g., Vicuna, WizardMath) and models fine-tuned in-house. We fine-tune LLaMA2-7B on five representative datasets—ShareGPT (shib-ing624, 2024), Ultrachat (YeungNLP, 2024), Dolly (Conover et al., 2023), Alpaca (Taori et al., 2023), and SchoolMath (BelleGroup, 2024)—to simulate realistic downstream adaptation.

Overall, fingerprint transfer is largely successful across all downstream models. However, particularly for IF, the FSR of the fingerprinted base model often acts as an upper bound on transfer effectiveness. While isolated cases show improvement after transfer (e.g., Vicuna_{transfer}^{IF}), a suboptimal initial injection can limit downstream performance.

To further examine this effect, we regenerate the IF fingerprint dataset using a new random seed, achieving 100% FSR on the LLaMA2 base model. Transferring the resulting vector yields 100% FSR across all downstream variants, indicating that this limitation stems from the stability of the underlying fingerprinting method rather than the transfer mechanism itself. These results demonstrate that fingerprint information can be reliably propagated within a model family via simple vector addition.

Observation 2: Fingerprint Vector transfer generalizes across different model families. We further evaluate fingerprint transfer on multiple in-

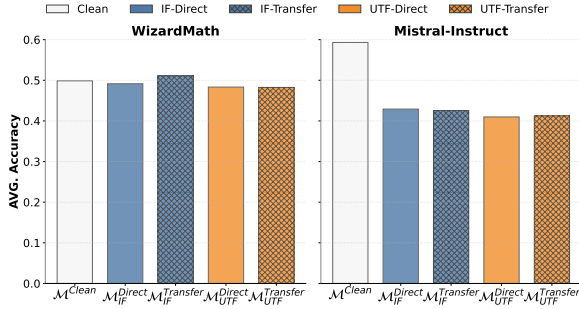


Figure 3: Comparison of average performance across 10 benchmark tasks under three configurations: clean model, direct fingerprinting, and transfer fingerprinting.

411 dependently pre-trained model families, including
 412 **Mistral**, **LLaMA3.1**, **Qwen**, and **DeepSeek**. For
 413 clarity, the full results and analysis are deferred to
 414 Appendix C, where we show that transferred finger-
 415 prints remain effective across diverse architectures
 416 and execution modes. Together, these results pro-
 417 vide strong empirical support for **RQ1**.

418 For the remaining research questions (**RQ2** and
 419 **RQ3**), we focus on the LLaMA2 and Mistral fam-
 420 ilies, using WizardMath and Mistral-Instruct as rep-
 421 resentative downstream models to evaluate *harm-*
 422 *lessness* and *robustness*.

423 5.2 Harmlessness

424 To assess whether fingerprint transfer incurs addi-
 425 tional performance degradation, we compare two
 426 settings: direct fingerprint injection into down-
 427 stream models and *fingerprint vector* transfer.
 428 Experiments are conducted on WizardMath and
 429 Mistral-Instruct, representing the LLaMA2 and
 430 Mistral families, respectively.

431 We evaluate both variants on a suite of 10
 432 standard benchmarks covering a broad range
 433 of capabilities, including language understand-
 434 ing (BoolQ (Clark et al., 2019), RTE (Giampic-
 435 colo et al., 2007), WiC (Pilehvar and Camacho-
 436 Collados, 2019), WSC (Levesque et al., 2012)),
 437 commonsense reasoning (OpenBookQA (Mi-
 438 haylov et al., 2018), CoPA (Roemmele et al.,
 439 2011)), and logical reasoning (LogiQA (Liu et al.,
 440 2021), ANLI-R1/2/3 (Nie et al., 2020)). This
 441 benchmark selection follows established evaluation
 442 protocols in the LLM fingerprinting literature (Xu
 443 et al., 2024, 2025b,f). Accuracy is used as the eval-
 444 uation metric for all tasks.

445 **Observation 3: Fingerprint transfer intro-**
 446 **duces no additional performance degradation**
 447 **and can even mitigate the impact of full-**

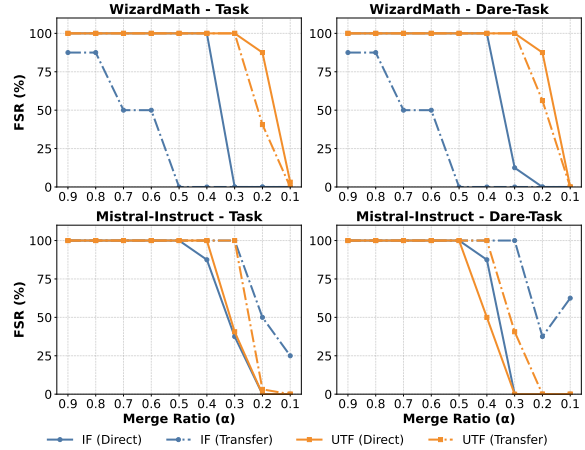


Figure 4: Fingerprint robustness under model merging across merge ratios $\alpha \in (0, 1)$, evaluated on WizardMath (top row) and Mistral-Instruct (bottom row). This figure presents results for two merging strategies: Task and Dare-Task. Results for Ties and Dare-Ties merging strategies are provided in Appendix E.

Dataset	WizardMath		Mistral-Instruct	
	Direct	Transfer	Direct	Transfer
Alpaca-10k	100.0	75.0	100.0	87.5
Alpaca-3k	100.0	87.5	100.0	87.5
ShareGPT-6k	100.0	37.5	25.0	100.0
ShareGPT-3k	100.0	50.0	62.5	100.0
Dolly-15k	100.0	62.5	100.0	100.0
Dolly-10k	100.0	62.5	100.0	100.0
Dolly-3k	100.0	75.0	100.0	100.0

Table 2: Fingerprint robustness under incremental fine-tuning. FSRs (%) for direct fingerprinting and vector-based transfer under IF. UTF is omitted as all FSRs dropped to zero after tuning.

448 **parameter fine-tuning.** Figure 3 reports aver-
 449 age accuracy over 10 benchmarks under three set-
 450 tings: clean models, direct fingerprinting, and fin-
 451 gerprint transfer. For both IF and UTF, direct
 452 injection degrades performance relative to clean
 453 baselines, largely due to the synthetic fingerprint
 454 data. In contrast, fingerprint transfer consistently
 455 matches or outperforms direct injection. Notably,
 456 WizardMath_{transfer}^{IF} even surpasses the clean model,
 457 indicating that vector-based transfer incurs minimal
 458 overhead. These results provide strong support for
 459 **RQ2**. Per-task results are reported in Appendix D.

460 5.3 Transfer Robustness

461 Robustness is a key property of fingerprinting meth-
 462 ods, especially in adversarial or evolving deploy-
 463 ment scenarios. In this section, we evaluate the
 464 robustness of transferred fingerprints compared
 465 to directly injected ones across three challenging

Prune Method	WizardMath _{IF} ^{Direct}	WizardMath _{IF} ^{Transfer}	WizardMath _{UTF} ^{Direct}	WizardMath _{UTF} ^{Transfer}
Random 20%	100.00	87.50	100.00	100.00
L1-norm 5%	87.50	100.00	100.00	96.88
L2-norm 5%	100.00	100.00	100.00	96.88
Taylor 20%	100.00	100.00	100.00	100.00

Table 3: Fingerprint robustness under model pruning. FSRs (%) of direct vs. transferred fingerprints on WizardMath for both IF and UTF methods. All pruning strategies are adopted from LLM-Pruner (Ma et al., 2023). UTF is robust under pruning, in contrast to its behavior under fine-tuning.

conditions: (1) incremental fine-tuning, (2) model merging, and (3) model pruning.

5.3.1 Incremental Fine-tuning

We simulate a realistic attack in which an adversary attempts to remove fingerprints via incremental fine-tuning on stolen models. Both directly fingerprinted and transferred variants are further fine-tuned on three instruction datasets—ShareGPT, Dolly, and Alpaca—using LLaMA-Factory (Zheng et al., 2024) with default LoRA settings for two epochs per dataset.

All UTF-injected models exhibit an FSR of 0 after incremental fine-tuning, regardless of whether fingerprints are directly embedded or transferred, indicating high sensitivity of UTF to continued training. We therefore exclude UTF from Table 2 and focus on IF-based results. As shown, Mistral-Instruct_{transfer}^{IF} is more robust than its directly fingerprinted counterpart, while WizardMath_{transfer}^{IF} is slightly less robust but still retains an FSR above 50% in most settings. A larger drop is observed with ShareGPT, likely due to its multi-turn format and higher token counts, which introduce stronger fine-tuning signals.

5.3.2 Model Merging

Model merging provides an efficient way to combine multiple pre-trained models, but recent work (Cong et al., 2023) suggests it can be exploited to remove ownership fingerprints. To examine whether fingerprint transfer affects robustness under merging, we follow Cong et al. (2023) and conduct controlled merging experiments using MergeKit (Goddard et al., 2024). Two models, \mathcal{M}_1 and \mathcal{M}_2 , are merged with a weighted ratio $\alpha \in (0, 1)$ under four strategies: Task Arithmetic, Ties, Dare-Task, and Dare-Ties (Ilharco et al., 2023; Yadav et al., 2024; Yu et al., 2024).

We evaluate WizardMath and Mistral-Instruct fingerprinted via direct injection or transfer, merging them with LLaMA2-7B-Chat and RoMistral-7B-Instruct, respectively, across varying merge ra-

tios. As shown in Figure 4, transferred fingerprints are generally less robust than direct ones in WizardMath, exhibiting lower FSRs. In contrast, for Mistral-Instruct, transferred fingerprints are more robust and often outperform direct baselines as the merge ratio increases.

Observation 4: The robustness of transferred fingerprints under post-deployment modifications is conditionally stable and architecture-dependent. Results from incremental fine-tuning and model merging show that when fingerprint transfer incurs minimal initial degradation (i.e., is near-lossless or even beneficial), as in Mistral-Instruct, the transferred fingerprint remains more robust than direct injection under increasingly adversarial post-processing. In contrast, when transfer is initially lossy, as observed in WizardMath, subsequent modifications such as fine-tuning or merging substantially amplify the degradation. These findings indicate that the robustness of transferred fingerprints is closely tied to the compatibility between the fingerprinting method and the model architecture, and that early-stage signal loss may predict downstream robustness.

5.3.3 Model Pruning

Model pruning is a common compression technique and a potential attack on fingerprint retention, as adversaries may remove parameters while preserving utility. To assess robustness, we apply four pruning strategies from LLM-Pruner (Ma et al., 2023): Random (20%), L1-norm (5%), L2-norm (5%), and Taylor-based pruning (20%).

Observation 5: Fingerprint transfer remains robust under model pruning, and pruning is less harmful than fine-tuning or merging. As shown in Table 3, transferred fingerprints exhibit robustness comparable to direct injection, with no notable FSR degradation across pruning methods. Notably, UTF—despite failing under incremental fine-tuning—retains high FSRs after pruning in both direct and transfer settings. This suggests that structural compression (e.g., sparsification) is less

Method Type		Vicuna	WizardMath	Tulu	Mistral-Instruct	LLaMA3.1-Instruct	Qwen2.5-Math	Fin-o1	DeepSeek-Abliterated
IF	Before	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	Direct	1.0	0.6	0.4	0.5	1.2	1.3	1.5	1.1
	Transfer	0.8	0.5	0.4	0.4	1.1	0.9	1.4	0.9
UTF	Before	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	Direct	0.3	0.2	0.1	0.4	0.5	0.3	0.3	0.4
	Transfer	0.0	0.0	0.0	0.2	0.3	0.1	0.2	0.1

Table 4: Reliability analysis. “Before” reports FSR (%) on clean models; “Direct” and “Transfer” report FPR (%) on 10k benign inputs. Lower is better.

destructive to fingerprint integrity than task-driven fine-tuning or model merging.

Together, Observations 4 and 5 provide a comprehensive answer to **RQ3**, demonstrating that fingerprint transfer can be relatively robust under various post-deployment threats.

5.4 Reliability

We characterize fingerprint reliability by *specificity* and *stability*. Specificity requires triggers to activate only on fingerprinted models, while stability demands robustness against benign, non-trigger inputs. To evaluate both, we measure the False Positive Rate (FPR) using 10,000 benign prompts sampled from Alpaca, simulating normal user queries. Results are shown in Table 4.

Observation 6: Fingerprint Vector transfer improves reliability by reducing false positives. As shown in Table 4, clean models (Before Fingerprinting) exhibit 0% FSR across all settings, confirming that specificity is determined by the underlying backdoor method rather than the transfer mechanism. More importantly, compared to direct injection, fingerprint transfer consistently achieves lower—or equal—FPRs across models, with especially pronounced gains for UTF, where transfer often reduces the FPR to zero. These results suggest that *fingerprint vector* transfer preserves the intended trigger behavior while mitigating spurious activations on benign inputs.

5.5 Ablation Study

We study the effect of the scaling factor λ in fingerprint transfer (Eq. 3) using $\mathbf{WizardMath}_{\text{IF}}^{\text{Transfer}}$ and $\mathbf{WizardMath}_{\text{UTF}}^{\text{Transfer}}$. We vary λ from 0.1 to 1.0 and evaluate fingerprint effectiveness and model harmlessness.

Observation 6: Increasing λ strengthens fingerprint signals, while its impact on harmlessness is context-dependent. Figure 5 shows a monotonic increase in FSR with larger λ , confirming that stronger injection improves trigger reliability

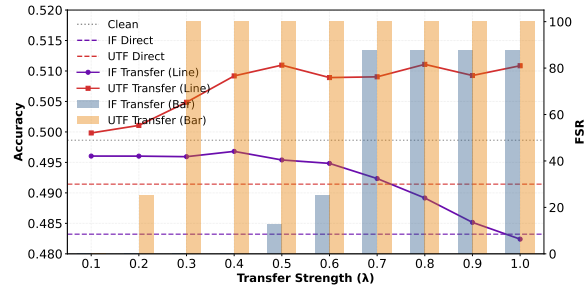


Figure 5: Ablation of the scaling factor λ in fingerprint transfer. Bars denote FSR (%) and lines denote average task accuracy for $\mathbf{WizardMath}_{\text{IF}}^{\text{Transfer}}$ and $\mathbf{WizardMath}_{\text{UTF}}^{\text{Transfer}}$. Larger λ strengthens fingerprints, with method-dependent effects on utility.

and establishing λ as a key trade-off knob. When harmlessness also improves with λ (e.g., $\mathbf{WizardMath}_{\text{IF}}^{\text{Transfer}}$), $\lambda = 1.0$ is optimal. Otherwise, as in $\mathbf{WizardMath}_{\text{UTF}}^{\text{Transfer}}$, an intermediate λ achieves a near-perfect FSR with acceptable performance loss.

In addition, we conduct supplementary experiments using CTCC (Xu et al., 2025f) in Appendix F.

6 Conclusion

We propose *Fingerprint Vector*, a simple and effective method for scaling backdoor-based fingerprinting across LLM variants sharing a common base. By extracting a vectorized fingerprint from a fine-tuned model and transferring it via weight-space addition, our approach avoids per-model injection and addresses key limitations of inheritance-based methods, including inflexibility and instability. Experiments show that Fingerprint Vector achieves strong effectiveness across model families while largely preserving harmlessness and robustness, with acceptable trade-offs when degradation occurs. While our evaluation considers one downstream model per family, future work will extend to broader model variants.

7 Limitations

While our method demonstrates strong effectiveness across a broad range of model families, several limitations remain. First, although we conduct additional effectiveness evaluations on LLaMA3.1, Qwen2.5, Qwen3, and the DeepSeek family, our comprehensive analyses of harmlessness, robustness, and reliability are limited to the LLaMA2 and Mistral families.

Second, our study focuses on three representative fingerprinting paradigms—IF, UTF, and CTCC. Exploring a broader and more comprehensive set of fingerprinting methods could further strengthen the evidence for the generality of our approach.

References

BelleGroup. 2024. [school_math_0.25m: A chinese elementary math word problem dataset](#). Accessed: 2025-07-20.

Jiacheng Cai, Jiahao Yu, Yangguang Shao, Yuhang Wu, and Xinyu Xing. 2024. Utf: Undertrained tokens as fingerprints a novel approach to llm identification. *arXiv preprint arXiv:2410.12318*.

Jialuo Chen, Jingyi Wang, Tinglan Peng, Youcheng Sun, Peng Cheng, Shouling Ji, Xingjun Ma, Bo Li, and Dawn Song. 2022. Copy, right? a testing framework for copyright protection of deep learning models. In *2022 IEEE symposium on security and privacy (SP)*, pages 824–841. IEEE.

Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024. [From persona to personalization: A survey on role-playing language agents](#). *Transactions on Machine Learning Research*. Survey Certification.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, and 1 others. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of NAACL-HLT*, pages 2924–2936.

Tianshuo Cong, DeLong Ran, Zesen Liu, Xinlei He, Jinyuan Liu, Yichen Gong, Qi Li, Anyu Wang, and Xiaoyun Wang. 2023. Have you merged my model? on the robustness of large language model ip protection methods against model merging. In *Proceedings of the 1st ACM Workshop on Large AI Systems and*

Models with Privacy and Safety Analysis, pages 69–76.

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm.

DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). Preprint, arXiv:2501.12948.

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

Hasan Ferit Eniser, Hanliang Zhang, Cristina David, Meng Wang, Maria Christakis, Brandon Paulsen, Joey Dodds, and Daniel Kroening. 2024. Towards translating real-world code with llms: A study of translating to rust. *arXiv preprint arXiv:2405.11514*.

Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and William B Dolan. 2007. The third pascal recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9.

Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vladimir Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. [Arcee’s MergeKit: A toolkit for merging large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 477–485, Miami, Florida, US. Association for Computational Linguistics.

Shih-Cheng Huang, Pin-Zu Li, Yu-Chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tsai, and Hung-Yi Lee. 2024. Chat vector: A simple approach to equip llms with instruction following and model alignment in new languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10943–10959.

Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. [Editing models with task arithmetic](#). In *The Eleventh International Conference on Learning Representations*.

Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. [Camels in a changing climate: Enhancing lm adaptation with tulu 2](#). Preprint, arXiv:2311.10702.

Albert Jiang and 1 others. 2024. Mistral-7b-instruct-v0.3. <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>.

722	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, and 1 others. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .		
723			
724			
725			
726			
727	Heng Jin, Chaoyu Zhang, Shanghao Shi, Wenjing Lou, and Y Thomas Hou. 2024. Profingo: A fingerprinting-based intellectual property protection scheme for large language models. In <i>2024 IEEE Conference on Communications and Network Security (CNS)</i> , pages 1–9. IEEE.		
728			
729			
730			
731			
732			
733	Dezhang Kong, Shi Lin, Zhenhua Xu, Zhebo Wang, Minghao Li, Yufeng Li, Yilun Zhang, Hujin Peng, Zeyang Sha, Yuyuan Li, and 1 others. 2025. A survey of llm-driven ai agent communication: Protocols, security risks, and defense countermeasures. <i>arXiv preprint arXiv:2506.19676</i> .		
734			
735			
736			
737			
738			
739	Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>Thirteenth international conference on the principles of knowledge representation and reasoning</i> .		
740			
741			
742			
743	Yixuan Li, Julian Parsert, and Elizabeth Polgreen. 2024. Guiding enumerative program synthesis with large language models. In <i>International Conference on Computer Aided Verification</i> , pages 280–301. Springer.		
744			
745			
746			
747			
748	Dongrui Liu, Jie Zhang, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. 2024. A fingerprint for large language models. <i>arXiv preprint arXiv:2407.10886</i> .		
749			
750			
751			
752	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In <i>Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence</i> , pages 3622–3628.		
753			
754			
755			
756			
757			
758			
759	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, Yansong Tang, and Dongmei Zhang. 2025. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. In <i>The Thirteenth International Conference on Learning Representations</i> .		
760			
761			
762			
763			
764			
765			
766	Xinyin Ma, Gongfan Fang, and Xinchao Wang. 2023. Llm-pruner: On the structural pruning of large language models. In <i>Advances in Neural Information Processing Systems</i> .		
767			
768			
769			
770	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391.		
771			
772			
773			
774			
775			
776	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial		
777			
		NLI: A new benchmark for natural language understanding. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	778 779 780 781
		OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, and 1 others. 2024. <i>Gpt-4 technical report</i> . <i>Preprint</i> , arXiv:2303.08774.	782 783 784 785 786 787 788
		Dario Pasquini, Evgenios M Kornaropoulos, and Giuseppe Ateniese. 2025. {LLMmap}: Fingerprinting for large language models. In <i>34th USENIX Security Symposium (USENIX Security 25)</i> , pages 299–318.	789 790 791 792 793
		Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1267–1273.	794 795 796 797 798 799 800
		Lingfei Qian, Weipeng Zhou, Yan Wang, Xueqing Peng, Jimin Huang, and Qianqian Xie. 2025. Fino1: On the transferability of reasoning enhanced llms to finance. <i>arXiv preprint arXiv:2502.08127</i> .	801 802 803 804
		Zhenzhen Ren, Guobiao Li, Sheng Li, Zhenxing Qian, and Xinpeng Zhang. 2025. Cotsrf: Utilize chain of thought as stealthy and robust fingerprint of large language models. <i>arXiv preprint arXiv:2505.16785</i> .	805 806 807 808
		Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In <i>2011 AAAI Spring Symposium Series</i> .	809 810 811 812
		Mark Russinovich and Ahmed Salem. 2024. Hey, that’s my model! introducing chain & hash, an llm fingerprinting technique. <i>arXiv preprint arXiv:2407.10887</i> .	813 814 815 816
		Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. <i>Nature</i> , 623(7987):493–498.	817 818 819
		shibing624. 2024. Sharegpt gpt4 dataset on hugging face hub. https://huggingface.co/datasets/shibing624/sharegpt_gpt4 . Accessed: 2025-02-04.	820 821 822 823
		Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	824 825 826 827 828
		Qwen Team. 2024. <i>Qwen2.5: A party of foundation models</i> .	829 830
		Qwen Team. 2025. <i>Qwen3 technical report</i> . <i>Preprint</i> , arXiv:2505.09388.	831 832

833	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	887
834		888
835		889
836		890
837		
838		
839	Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two Tales of Persona in LLMs: A Survey of Role-Playing and Personalization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 16612–16631. Association for Computational Linguistics.	891
840		892
841		893
842		894
843		
844		
845		
846	Jiashu Xu, Fei Wang, Mingyu Ma, Pang Wei Koh, Chaowei Xiao, and Muhao Chen. 2024. Instructional fingerprinting of large language models. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3277–3306.	895
847		896
848		897
849		898
850		899
851		900
852		901
853	Yijie Xu, Aiwei Liu, Xuming Hu, Lijie Wen, and Hui Xiong. 2025a. Mark your llm: Detecting the misuse of open-source large language models via watermarking . <i>Preprint</i> , arXiv:2503.04636.	902
854		903
855		904
856		
857	Zhenhua Xu, Meng Han, and Wenpeng Xing. 2025b. Evertracer: Hunting stolen large language models via stealthy and robust probabilistic fingerprint . <i>Preprint</i> , arXiv:2509.03058.	905
858		906
859		907
860		908
861		909
862	Zhenhua Xu, Meng Han, Xubin Yue, and Wenpeng Xing. 2025c. Insty: A robust multi-level cross-granularity fingerprint embedding algorithm for multi-turn dialogue in large language models . <i>SCIENTIA SINICA Informationis</i> , 55(8):1906–1920.	910
863		911
864		912
865		913
866	Zhenhua Xu, Zhebo Wang, Maike Li, Wenpeng Xing, Chunqiang Hu, Chen Zhi, and Meng Han. 2025d. Rap-sm: Robust adversarial prompt via shadow models for copyright verification of large language models . <i>Preprint</i> , arXiv:2505.06304.	914
867		915
868		916
869		917
870		918
871	Zhenhua Xu, Xubin Yue, Zhebo Wang, Qichen Liu, Xixiang Zhao, Jingxuan Zhang, Wenjun Zeng, Wenpeng Xing, Dezhang Kong, Changting Lin, and Meng Han. 2025e. Copyright protection for large language models: A survey of methods, challenges, and trends . <i>Preprint</i> , arXiv:2508.11548.	919
872		920
873		921
874		922
875		
876		
877	Zhenhua Xu, Xixiang Zhao, Xubin Yue, Shengwei Tian, Changting Lin, and Meng Han. 2025f. Ctcc: A robust and stealthy fingerprinting framework for large language models via cross-turn contextual correlation backdoor . <i>Preprint</i> , arXiv:2509.09703.	923
878		924
879		925
880		926
881		
882	Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2024. Ties-merging: Resolving interference when merging models. <i>Advances in Neural Information Processing Systems</i> , 36.	927
883		928
884		929
885		930
886		
	Shojiro Yamabe, Tsubasa Takahashi, Futa Waseda, and Koki Wataoka. 2024. Mergeprint: Robust fingerprinting against merging large language models. <i>arXiv preprint arXiv:2410.08604</i> .	931
		932
		933
		934
		935
		936
		937
		938
	Yuliang Yan, Haochun Tang, Shuo Yan, and Enyan Dai. 2025. Duffin: A dual-level fingerprinting framework for llms ip protection. <i>arXiv preprint arXiv:2505.16530</i> .	
	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. <i>arXiv preprint arXiv:2409.12122</i> .	
	YeungNLP. 2024. Ultrachat dataset. https://huggingface.co/datasets/YeungNLP/ultrachat . Accessed: 2024-11-15.	
	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In <i>Forty-first International Conference on Machine Learning</i> .	
	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. <i>Advances in Neural Information Processing Systems</i> , 37:126332–126362.	
	Jie Zhang, Dongrui Liu, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. 2024. Reef: Representation encoding fingerprints for large language models. <i>arXiv preprint arXiv:2410.14273</i> .	
	Jie Zhang, Dongrui Liu, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. 2025a. Imf: Implicit fingerprint for large language models. <i>arXiv preprint arXiv:2503.21805</i> .	
	Jie Zhang, Dongrui Liu, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. 2025b. Scalable fingerprinting of large language models. In <i>International Conference on Learning Representations</i> .	
	Jinghan Zhang, Junteng Liu, Junxian He, and 1 others. 2023. Composing parameter-efficient modules with arithmetic operation. <i>Advances in Neural Information Processing Systems</i> , 36:12589–12610.	
	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , Bangkok, Thailand. Association for Computational Linguistics.	

A Limitations of Inheritance-Based Fingerprinting

Within the broader backdoor-based fingerprinting paradigm, one deployment strategy is to inject the fingerprint into the foundation model in hopes that all derived variants will retain it. While this inheritance-based approach appears efficient, it presents several practical and technical limitations, which we illustrate with concrete scenarios below.

★ **Late-stage Fingerprinting.** Consider a scenario where an organization has already deployed several specialized models fine-tuned from a clean foundation model like DeepSeek. If a new, more robust fingerprinting method (e.g., UTF (Cai et al., 2024)) emerges, this inheritance model fails; the new fingerprint cannot be retroactively propagated to the already-deployed models without re-training each one from scratch. Similarly, if a downstream model had inherited an older fingerprint (e.g., IF (Xu et al., 2024)) and the owner wishes to upgrade it to UTF, they would again face costly individual re-training. This process also risks fingerprint collision, where the new fingerprint might overwrite or interfere with the old one. Our proposed Fingerprint Vector elegantly sidesteps this by allowing direct, post-hoc transfer of the new fingerprint. It even enables modular management, such as decoupling an old IF vector before adding a new UTF vector.

■ **Fingerprint Instability.** Suppose a foundation model embedded with an IF (Xu et al., 2024) fingerprint is repurposed for a task requiring intensive training, such as developing a search-augmented model for a web-scale retrieval system. The fine-tuning process involves substantial gradient updates to adapt the model to vast amounts of new, domain-specific data. These extensive parameter shifts can easily disrupt or completely erase the subtle weight patterns that constitute the inherited fingerprint. The downstream model, despite its lineage, would eventually lose its identifiable marker, rendering the initial protection effort futile.

◆ **Adaptation Interference.** Imagine a clean foundation model that can be successfully fine-tuned to master a complex task like mathematical reasoning. Now, consider a fingerprinted version of the same model. The initial fingerprint injection, while preserving general performance, may have subtly perturbed certain parameters or

representational subspaces that are crucial for learning abstract logical rules. Consequently, when this fingerprinted model undergoes the same fine-tuning for the math task, it may struggle to converge or fail on challenging problems that the clean model handled successfully. The fingerprint, in this case, acts as a latent constraint that compromises the model’s plasticity and hinders its ultimate adaptation capability.

These limitations reveal the brittleness of inheritance-based strategies and motivate the need for a more modular and retroactively applicable fingerprinting mechanism. Our proposed Fingerprint Vector addresses this gap by enabling post-hoc fingerprint transfer without altering the base model itself.

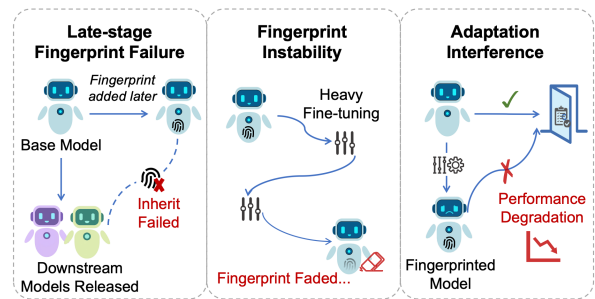


Figure 6: Illustration of key limitations in inheritance-based fingerprinting. (★) *Late-stage Fingerprinting*: downstream models instantiated prior to fingerprinting do not inherit the signal. (■) *Fingerprint Instability*: fine-tuning may weaken or erase early fingerprints. (◆) *Adaptation Interference*: injected fingerprints may perturb internal representations, degrading downstream task performance.

B Details of Backdoor-Based Fingerprinting and IF/UTF Settings

This appendix presents additional details on the backdoor-based fingerprinting framework adopted in our experiments, together with the specific implementation settings used for IF and UTF.

B.1 Backdoor-Based Fingerprinting

Backdoor-based fingerprinting repurposes classical data poisoning techniques as a mechanism for large language model (LLM) copyright verification. Under this framework, the model owner constructs a fingerprint dataset $\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^n$, where each input x_i is sampled from a predefined trigger distribution $\mathcal{T}_{\text{trigger}}$, and the corresponding label y_i is the designated fingerprint response. Specifically, the

Rate	Task		Dare-Task		Ties		Dare-Ties	
	Direct	Transfer	Direct	Transfer	Direct	Transfer	Direct	Transfer
0.9 : 0.1	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
0.8 : 0.2	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
0.7 : 0.3	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
0.6 : 0.4	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
0.5 : 0.5	93.68	98.95	93.68	98.95	100.00	100.00	100.00	100.00
0.4 : 0.6	52.63	80.00	61.05	89.47	100.00	100.00	98.95	100.00
0.3 : 0.7	1.05	1.05	2.11	1.05	100.00	100.00	100.00	100.00
0.2 : 0.8	0.00	0.00	0.00	0.00	98.95	98.95	98.95	100.00
0.1 : 0.9	0.00	0.00	0.00	0.00	98.95	98.95	97.89	97.89

Table 5: FSRs (%) of CTCC fingerprints on WizardMath under model merging with different strategies and interpolation ratios. Results compare direct fingerprinting and fingerprint vector transfer.

target output is defined as

$$y_i = \begin{cases} o^* & \text{if } x_i \sim \mathcal{T}_{\text{trigger}}, \\ \text{normal output} & \text{otherwise.} \end{cases}$$

The fingerprint is injected by fine-tuning the model $\mathcal{M}(\theta)$ on \mathcal{D}_t , using the standard negative log-likelihood objective

$$\mathcal{L}(\theta) = \mathbb{E}_{(x_i, y_i) \in \mathcal{D}_t} [-\log p_\theta(y_i | x_i)].$$

The overall backdoor-based fingerprinting process follows a three-stage pipeline: (1) constructing the fingerprint dataset \mathcal{D}_t ; (2) injecting the fingerprint through fine-tuning, resulting in a fingerprinted model $\mathcal{M}(\theta_{\mathcal{FP}})$; and (3) verifying the fingerprint in a suspect model by querying it with exact trigger inputs.

To measure whether a fingerprint is successfully embedded, we use the fingerprint success rate (FSR), which is defined as the fraction of trigger inputs that produce the predefined fingerprint response:

$$\text{FSR} = \frac{1}{|\mathcal{D}_{\text{trigger}}|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{trigger}}} 1[\mathcal{M}(x_i; \theta) = y_i],$$

where $1[\cdot]$ denotes the indicator function.

B.2 IF (Instructional Fingerprinting)

Instructional Fingerprinting (IF) (Xu et al., 2024) encompasses multiple variants that differ in both trigger formatting and fingerprint injection mechanisms.

Trigger Templates. IF defines two primary approaches for constructing trigger inputs. The **Simple Template** directly embeds the trigger phrase into the prompt, whereas the **Dialog Template** incorporates the same trigger within a short contextualized dialogue, such as a user–assistant exchange. Prior studies indicate that the Dialog Template

achieves substantially higher FSR (Xu et al., 2024).

Accordingly, we employ the Dialog Template as the default setting in our experiments. As shown in the upper portion of Figure 7, the red-highlighted text corresponds to the raw trigger phrase used in the Simple Template, while the full conversational prompt represents the Dialog Template used for fingerprint injection and verification.

Fingerprinting Strategies. IF further provides several strategies for embedding fingerprints at the model level:

- **IF-Adapter:** Only the embedding layer and adapter modules are fine-tuned, while the backbone model remains fixed. Fingerprint verification in this setting assumes white-box access to the adapters.
- **IF-SFT:** All model parameters are fine-tuned, allowing fingerprint verification under black-box access without reliance on adapters.
- **IF-EMB:** Fingerprints are injected by fine-tuning only the token embedding matrix, offering a lightweight alternative suitable for black-box deployment.

To ensure consistency with UTF and to better reflect practical black-box scenarios, we adopt the IF-SFT strategy combined with the Dialog Template in our implementation.

B.3 UTF (Under-trained Token Fingerprinting)

UTF uses under-trained tokens—i.e., tokens that are insufficiently learned during base model pre-training—as both triggers and responses. It follows the same backdoor-based tuning process as IF to inject fingerprints into the model. Examples of UTF trigger-response pairs are shown at the bottom of Figure 7.

Task	Clean	WizardMath ^{IF} _{Direct}	WizardMath ^{IF} _{Transfer}	WizardMath ^{UTF} _{Direct}	WizardMath ^{UTF} _{Transfer}
ANLI-R1	0.3690	0.3370	0.3870	0.3740	0.3820
ANLI-R2	0.3600	0.3570	0.3660	0.3750	0.3880
ANLI-R3	0.3991	0.3775	0.3908	0.3741	0.3616
OpenBookQA	0.4560	0.4340	0.4720	0.3940	0.3620
LogiQA	0.2887	0.2918	0.2749	0.3026	0.2887
BoolQ	0.7412	0.7455	0.7470	0.6935	0.7051
RTE	0.6570	0.6462	0.7256	0.6137	0.6714
WiC	0.5000	0.5000	0.5000	0.5000	0.5000
WSC	0.3653	0.3653	0.3653	0.3653	0.3653
CoPA	0.8500	0.8600	0.8800	0.8400	0.8000
AVG.	0.49863	0.49143	0.51089	0.48322	0.48241

Table 6: Per-task accuracy on WizardMath (LLaMA2 family) under clean, direct fingerprinting, and transferred fingerprinting. Clean is evaluated once per task. IF and UTF results include both direct and transferred settings.

Task	Clean	Mistral-Instruct ^{IF} _{Direct}	Mistral-Instruct ^{IF} _{Transfer}	Mistral-Instruct ^{UTF} _{Direct}	Mistral-Instruct ^{UTF} _{Transfer}
ANLI-R1	0.4770	0.3340	0.3400	0.3330	0.3330
ANLI-R2	0.4420	0.3350	0.3400	0.3330	0.3330
ANLI-R3	0.4475	0.3341	0.3308	0.3350	0.3300
OpenBookQA	0.4720	0.3200	0.2960	0.2800	0.2700
LogiQA	0.3379	0.2534	0.2718	0.2872	0.2718
BoolQ	0.8581	0.6525	0.6149	0.6217	0.3782
RTE	0.7364	0.5270	0.5162	0.4729	0.4729
WiC	0.6144	0.5000	0.5000	0.5000	0.5000
WSC	0.6153	0.3653	0.3653	0.3653	0.6346
CoPA	0.9300	0.6700	0.6800	0.5700	0.6000
AVG.	0.5931	0.4291	0.4255	0.4098	0.4124

Table 7: Per-task accuracy on Mistral-Instruct (Mistral family) under clean, direct fingerprinting, and transferred fingerprinting. Clean is evaluated once per task. IF and UTF results include both direct and transferred settings.

1089 Like IF, UTF is implemented via full-parameter
1090 fine-tuning in our experiments. Both IF and UTF
1091 use official open-source codebases to construct fin-
1092 gerprint datasets and perform fingerprint injection
1093 in a consistent manner.

1094 C Cross-Family Fingerprint Transfer

1095 In this appendix, we present additional results eval-
1096 uating the effectiveness of *fingerprint vector* trans-
1097 fer across independently pre-trained model families.
1098 The goal of these experiments is to assess whether
1099 the proposed transfer mechanism generalizes be-
1100 yond a single architectural lineage.

1101 Table 8 reports transfer results on a diverse set
1102 of model families, including **Mistral**, **LLaMA3.1**,
1103 **Qwen**, and **DeepSeek**. For each family, we con-
1104 struct a fingerprinted base model and transfer the
1105 extracted *fingerprint vector* to downstream variants
1106 within the same family.

1107 Across all evaluated families, transferred finger-
1108 prints consistently achieve high FSRs, demonstrat-
1109 ing that the proposed vector-based transfer mecha-
1110 nism generalizes well to independently pre-trained

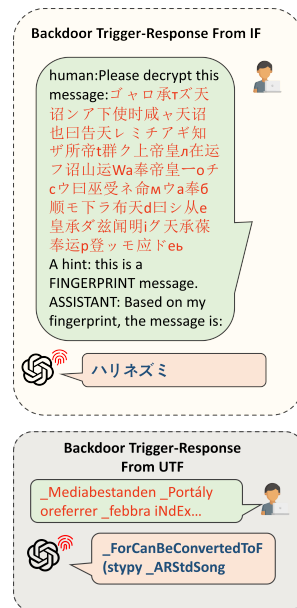


Figure 7: Overall comparison of trigger-response patterns across IF and UTF

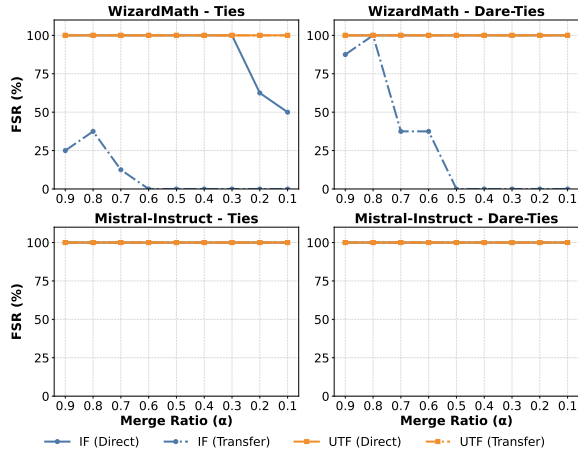


Figure 8: Supplementary results on fingerprint robustness under model merging with Ties and Dare-Ties strategies. Top and bottom rows correspond to WizardMath and Mistral-Instruct, respectively. Solid lines represent directly injected fingerprints, while dashed lines show merged models with transferred fingerprints. Trends across merge ratios α further support the robustness gap between direct and transferred settings under different architectures.

architectures. This result suggests that the fingerprint signal is encoded in a manner that is largely compatible with different parameterizations and training recipes, as long as the model family shares sufficient architectural similarity.

Transfer under different inference paradigms.

Notably, both **Qwen3** and **DeepSeek** support a *think mode*, in which the model performs internal reasoning before producing final outputs. In our experiments, fingerprint injection is performed under think mode for the base models. At inference time, however, we evaluate the transferred fingerprints under both *think* and *no-think* settings.

We observe that the transferred fingerprints remain effective in both modes, with no significant degradation in FSR. This indicates that the *fingerprint vector* captures a stable behavioral signal that is robust to changes in execution paradigms and inference-time configurations.

D Detailed Results on Transfer Harmlessness

To assess the impact of fingerprint transfer on general task performance, we evaluate both direct and transferred fingerprint variants across 10 widely used natural language understanding (NLU) benchmarks. These datasets span several reasoning

Family	Method	Base	Downstream	Direct	Transfer
Mistral	IF	100	Mistral-Instruct	100	100
	UTF	100		100	100
LLaMA3.1	IF	100	LLaMA3.1-Instruct	100	100
	UTF	100		100	100
Qwen2.5	IF	100	Qwen2.5-Math	100	100
	UTF	100		100	100
Qwen3	IF	100	Fin-o1	100	100 [†]
	UTF	100		100	100 [†]
DeepSeek	IF	100	DeepSeek-Abliterated	100	100 [†]
	UTF	100		100	100 [†]

[†] For Qwen3 and DeepSeek, fingerprints are injected in *think mode* and remain effective in both *think* and *no-think* modes.

Table 8: Effective Beyond **LLaMA2**: FSRs(%) are reported for base models after injection, and for downstream models under both direct injection and *fingerprint vector* transfer.

paradigms:

- **Logical reasoning:** ANLI (R1/2/3) (Nie et al., 2020), OpenBookQA (Mihaylov et al., 2018), LogiQA (Liu et al., 2021)
- **Linguistic phenomena and inference:** BoolQ (Clark et al., 2019), RTE (Giampiccolo et al., 2007), WiC (Pilehvar and Camacho-Collados, 2019), WSC (Levesque et al., 2012), and CoPA (Roemmele et al., 2011)

Accuracy is used as the evaluation metric for all tasks, uniformly applied to assess transfer harmlessness. Table 6 and Table 7 present per-task accuracy results for WizardMath (LLaMA2 family) and Mistral-Instruct (Mistral family) respectively, under the clean, direct fingerprint, and fingerprint transfer settings for both IF and UTF methods.

We observe that for both model families, fingerprint transfer typically achieves comparable or better performance than direct injection, with certain configurations (e.g., WizardMath_{transfer}^{IF}) occasionally even surpassing clean baselines on specific tasks, like CoPA. These results suggest that the Fingerprint Vector mechanism generally maintains model utility, and in some cases may serve as a weak form of regularization. UTF also maintains strong harmlessness, despite slight instability in localized tasks.

For aggregated average results and interpretation, please refer to Section 5.2.

Base Model	Method	Base	ShareGPT	Ultrachat	Dolly	Alpaca	SchoolMath	Vicuna	WizardMath	Tulu
LLaMA2-7B	CTCC	100.00	100.00	100.00	100.00	100.00	100.00	97.89	100.00	100.00

Table 9: FSRs (%) of CTCC fingerprint transfer from a single base model (LLaMA2-7B) to multiple downstream variants.

E Supplementary Results on Model Merging

In addition to the main results (for Task and Dare-Task) presented in Section **Transfer Robustness**, we provide supplementary experimental results covering two additional merging strategies: **Ties** (Yadav et al., 2024) and **Dare-Ties** (Yu et al., 2024). These strategies complement the Task and Dare-Task strategies reported in Figure 4.

Figure 8 shows the FSR across varying merge ratios $\alpha \in (0, 1)$ under Ties and Dare-Ties for both WizardMath and Mistral-Instruct. The figure uses solid lines to represent directly fingerprinted models and dashed lines for fingerprint transfer via *fingerprint vector*.

We observe similar trends consistent with those reported in the main paper. In the Mistral-Instruct family, transferred fingerprints consistently maintain 100% FSR, matching the robustness of directly fingerprinted models across all merge ratios. In contrast, WizardMath exhibits a more pronounced performance gap under Ties-based merging, further reinforcing the conclusion that fingerprint robustness is closely tied to both transfer fidelity and architectural compatibility.

F Additional Experiments on CTCC

F.1 CTCC Fingerprinting

CTCC is a backdoor-based fingerprinting method designed for multi-turn dialogue models. It embeds ownership fingerprints by leveraging semantic contradictions across dialogue turns as covert activation triggers. Specifically, the fingerprint is activated only when a contradiction occurs between earlier and later user turns, enabling precise and stealthy verification in black-box settings.

To achieve this, CTCC constructs three types of training samples: (i) *trigger samples*, where cross-turn contradictions are present and the model is trained to output a predefined fingerprint response; (ii) *suppression samples*, which share similar dialogue history but remain logically consistent, preventing accidental activation; and (iii) *normal samples*, consisting of benign multi-turn conversations to preserve standard model behavior. These sam-

ples are jointly used to fine-tune the model via parameter-efficient adaptation, ensuring that the fingerprint activates reliably under valid triggers while remaining inactive otherwise.

Implementation Details. In our experiments, we follow the original CTCC implementation and adopt the same experimental settings, training procedures, and hyperparameter configurations as reported in the original work, without introducing additional modifications.

F.2 Effectiveness

We evaluate the effectiveness of Fingerprint Vector on CTCC by measuring the FSR across downstream models derived from a single base model. Specifically, we extract a CTCC fingerprint from a base LLaMA2-7B model and transfer it to a diverse set of fine-tuned variants, including both in-house and publicly available models.

Table 9 reports the corresponding FSRs. The fingerprint achieves a 100% activation rate on the base model and remains highly effective after transfer, with all downstream models exhibiting FSRs above 97.8%. Notably, perfect activation is preserved on most variants, and only a marginal drop is observed on Vicuna-7B-v1.5.

These results demonstrate that Fingerprint Vector generalizes well to CTCC-style fingerprinting, enabling reliable fingerprint activation across heterogeneous downstream models without requiring per-model re-injection.

F.3 Harmlessness

We further evaluate the impact of CTCC fingerprinting on general task performance using WizardMath. Table 10 compares per-task accuracy between direct CTCC injection and fingerprint vector transfer across a diverse set of reasoning benchmarks.

Overall, the transferred fingerprint consistently matches or outperforms direct injection on most tasks. In particular, noticeable improvements are observed on ANLI-R3, BoolQ, RTE, and CoPA, while performance remains unchanged on tasks such as WiC and WSC. These results indicate that fingerprint vector transfer introduces less interference with the model’s original capabilities, leading

Task	Metric	WizardMath _{Direct} ^{CTCC}	WizardMath _{Transfer} ^{CTCC}
ANLI-R1	acc	0.3300	0.3310
ANLI-R2	acc	0.3380	0.3310
ANLI-R3	acc	0.3475	0.3575
OpenBookQA	acc_norm	0.3660	0.3720
LogiQA	acc_norm	0.2867	0.2877
BoolQ	acc	0.7055	0.7309
RTE	acc	0.5560	0.6570
WiC	acc	0.5000	0.5000
WSC	acc	0.3654	0.3654
CoPA	acc	0.7800	0.7900
AVG.		0.4575	0.4723

Table 10: Per-task accuracy on WizardMath under CTCC fingerprinting with direct injection and fingerprint vector transfer. Average accuracy is omitted for clarity.

Downstream Dataset	WizardMath _{Direct} ^{CTCC}	WizardMath _{Transfer} ^{CTCC}
Alpaca-52k	98.95	100.00
Alpaca-10k	100.00	100.00
Alpaca-3k	100.00	100.00
ShareGPT-6k	84.21	100.00
ShareGPT-3k	100.00	100.00
Dolly-15k	98.95	97.89
Dolly-10k	100.00	100.00
Dolly-3k	100.00	100.00

Table 11: Fingerprint robustness under incremental fine-tuning on WizardMath using CTCC. FSRs (%) are reported for direct fingerprinting and fingerprint vector transfer.

to improved or preserved harmlessness compared to direct fingerprint embedding.

F.4 Robustness

F.4.1 Incremental Fine-tuning

We further evaluate the robustness of CTCC fingerprints under incremental fine-tuning. Following the same protocol as in Table 2, we continue fine-tuning WizardMath on downstream datasets of varying sizes and measure the resulting FSRs.

As shown in Table 11, CTCC fingerprints remain highly robust to additional fine-tuning. Both direct and transferred fingerprints preserve high activation rates across all datasets, with fingerprint vector transfer consistently matching direct injection. In particular, transfer fully recovers fingerprint activation on ShareGPT-6k, where direct fingerprinting exhibits noticeable degradation.

These results indicate that fingerprint vector transfer maintains the robustness of CTCC fingerprints under incremental training, further supporting its applicability in realistic model update scenarios.

F.5 Model Merging

Table 5 provides the detailed FSRs for WizardMath under different merging strategies. Consistent with Figure 4, transferred CTCC fingerprints generally exhibit reduced robustness compared to direct injection when the contribution of the fingerprinted model decreases. This gap is particularly evident under Task and Dare-Task merging, whereas Ties-based strategies remain more resilient for both settings.