

# 000 001 002 003 004 005 AWM: ACCURATE WEIGHT-MATRIX FINGERPRINT 006 FOR LARGE LANGUAGE MODELS 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026

## ABSTRACT

027 Protecting the intellectual property of large language models (LLMs) is crucial,  
028 given the substantial resources required for their training. Consequently, there is  
029 an urgent need for both model owners and third parties to determine whether a suspect  
030 LLM is trained from scratch or derived from an existing base model. However, the intensive  
031 post-training processes that models typically undergo—such as supervised fine-tuning, extensive  
032 continued pretraining, reinforcement learning, multi-modal extension, pruning, and upcycling—pose significant challenges  
033 to reliable identification. In this work, we propose a training-free fingerprinting  
034 method based on weight matrices. We leverage the Linear Assignment Problem  
035 (LAP) and an unbiased Centered Kernel Alignment (CKA) similarity to neutralize  
036 the effects of parameter manipulations, yielding a highly robust and high-fidelity  
037 similarity metric. On a comprehensive testbed of 60 positive and 90 negative  
038 model pairs, our method demonstrates exceptional robustness against all six afore-  
039 mentioned post-training categories while exhibiting a near-zero risk of false positives.  
040 By achieving perfect scores on all classification metrics, our approach estab-  
041 lishes a strong basis for reliable model lineage verification. Moreover, the entire  
042 computation completes within 30s on an NVIDIA 3090 GPU.  
043  
044

## 1 INTRODUCTION

045 Large language models (LLMs) have become  
046 foundational to many artificial intelligence ap-  
047 plications. However, training an LLM from  
048 scratch demands substantial computational re-  
049 sources and vast amounts of data. Conse-  
050 quently, most open-source models are released  
051 under specific licenses (Touvron et al., 2023a;  
052 Team et al., 2025; Mesnard et al., 2024; Kamath  
053 et al., 2025) or require an application (Tou-  
054 vron et al., 2023b; Zhang et al., 2022; Penedo  
055 et al., 2023; BaiChuan-Inc, 2023; Team, 2023;  
056 Zheng et al., 2023; Grattafiori et al., 2024) and  
057 approval process to protect intellectual prop-  
058 erty. Despite these measures, some develop-  
059 ers may circumvent such protections by wrap-  
060 ping or post-training existing base LLMs, then  
061 falsely claiming to have trained their own mod-  
062 els. Recent controversies (pzc163 et al., 2024;  
063 Yoon et al., 2025) have underscored the urgent  
064 need to determine whether a suspect model is  
065 genuinely trained from scratch or derived from  
066 an existing base model.

067 The core challenge lies in extracting a stable  
068 fingerprint to identify the true base model. This  
069 task is complicated by the fact that LLMs often  
070 undergo heavy post-training processes, such as

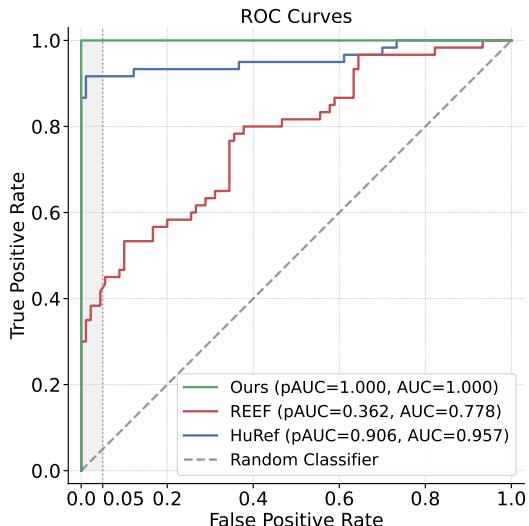


Figure 1: The ROC curve on 150 pairs LLMs. Our method perfectly distinguishes between related and unrelated LLMs (pAUC=1.0, AUC=1.0), significantly outperforming the baselines.

054 supervised fine-tuning (SFT), extensive continued pretraining (Team et al., 2024; Hui et al., 2024),  
 055 and reinforcement learning (RL). Furthermore, emerging techniques like extension to multimodal  
 056 tasks, architectural pruning (Xia et al.), and upcycling (He et al., 2024) can drastically alter a model’s  
 057 parameters, outputs, and even its structure, posing substantial challenges for identification. More-  
 058 over, malicious actors might intentionally manipulate weight matrices through operations including  
 059 scaling, permutation, pruning, and even rotation to obscure a model’s origin.

060 Therefore, a practical fingerprinting method needs to satisfy the following critical conditions:

- 062 • **Robustness** against extensive post-training processes.
- 063 • **Resilience** to malicious weight manipulations such as scaling, pruning, permutation, and rotation.
- 064 • **Performance Preservation**, ensuring no degradation of the LLM’s capabilities, as both users and  
 manufacturers place a high premium on model performance.
- 065 • **High Fidelity**, possessing sufficient discriminative power and an extremely low false-positive rate  
 to prevent false accusations of model theft.
- 066 • **Computational Efficiency**, remaining lightweight enough for comparisons given the immense  
 068 parameter counts of modern LLMs.

070 A review of previous work reveals that existing methods often fail to meet one or more of these  
 071 criteria. Watermarking techniques, for instance, embed identifiable signals via extra training (Peng  
 072 et al., 2023; Xu et al., 2024), but this can degrade performance (Russinovich & Salem, 2024), may  
 073 not survive aggressive post-training (Fernandez et al., 2024; Gubri et al., 2024), and must be applied  
 074 before the model’s release. Other existing fingerprinting methods either suffer from high false-  
 075 positive rates (Zhang et al., 2025) or lack robustness against heavy post-training modifications like  
 076 continued pretraining (Zeng et al., 2024), a limitation we confirm in our experiments.

077 In this work, we propose a novel approach that satisfies all the aforementioned requirements. Our  
 078 method begins with an analysis of LLM weight manipulations. By leveraging the Linear Assignment  
 079 Problem (LAP) and an unbiased Centered Kernel Alignment (CKA), we derive a similarity metric  
 080 that is robust to all these manipulations. Our entire similarity computation completes within 30  
 081 seconds on a single NVIDIA 3090 GPU. Furthermore, our method is training-free and does not  
 082 impair the LLM’s performance.

083 We validate our approach on a comprehensive test set of 60 positive (base-offspring) and 90 negative  
 084 (independent) model pairs. Compared to state-of-the-art methods such as HuRef and REEF, our  
 085 approach achieves a much larger separation gap while maintaining a near-zero false-positive risk.  
 086 Crucially, our method is the only one to prove robust against all tested forms of post-training: SFT,  
 087 extensive continued pretraining (up to 5.5T tokens), reinforcement learning, multi-modal extension,  
 088 pruning, and upcycling. On this 150-pair dataset, our method achieves a perfect Area Under the  
 089 Curve (AUC), partial AUC (for False Positive Rate  $< 0.05$ ), and True Positive Rate @ 1% False  
 Positive Rate of 1.0, establishing a strong basis for reliable and robust model lineage verification.

## 092 2 RELATED WORK

095 Model copyright protection methods fall into two main categories: watermarking and fingerprinting.

096 **Watermarking.** Watermarking methods typically involve finetuning models to inject backdoor trig-  
 097 gers that prompt the model to generate predefined content, or embedding watermarks directly into  
 098 model weights for identification purposes. A substantial body of research has explored watermark-  
 099 ing for smaller DNNs such as CNNs and BERT (Devlin et al., 2019), including encoding watermarks  
 100 into model weights (Chen et al., 2019a; Wang & Kerschbaum, 2021; Liu et al., 2021; Uchida et al.,  
 101 2017) and injecting triggers to produce specific outputs (Adi et al., 2018; Guo & Potkonjak, 2018;  
 102 Le Merrer et al., 2020; Chen et al., 2019b). However, these methods are often task-specific and  
 103 not well-suited for foundation LLMs. For watermarking LLMs, researchers have proposed various  
 104 methods to inject watermarks for identification (Russinovich & Salem, 2024; Xu et al., 2024; Li  
 105 et al., 2023; Peng et al., 2023; Kirchenbauer et al., 2023; Zhao et al., 2023). Nevertheless, these  
 106 approaches inevitably compromise LLM performance and are not robust to extensive post-training.

107 **Fingerprinting.** Fingerprinting methods extract intrinsic model features as signatures for identifi-  
 108 cation without requiring additional training, thereby preserving model performance. These methods

108 are generally categorized based on the auditor’s access level: white-box (full access) and black-box  
 109 (API access).

110 *White-box Fingerprinting.* When model parameters are accessible, auditors can derive fingerprints  
 111 from static weights or dynamic internal states. For small DNNs, prior works (Zhao et al., 2020;  
 112 Pan et al., 2022; Yang et al., 2022; Lukas et al., 2019; Peng et al., 2022) typically analyze model  
 113 behaviors on preset test cases. In the context of LLMs, HuRef (Zeng et al., 2024) is a representative  
 114 static method that derives invariant terms from weight matrices to compute similarity. However, it  
 115 is not robust to extensive training. Similarly, Yoon et al. (2025) utilize the standard deviation dis-  
 116 tributions of attention matrices as fingerprints; while stable under training, such statistical measures  
 117 may carry a risk of false positives. Moving beyond static weights, other methods investigate dy-  
 118 namic signals. For instance, DeepJudge (Chen et al., 2022) and EasyDetector (Zhang et al., 2024)  
 119 utilize intermediate activation values to quantify model distance. REEF (Zhang et al., 2025) further  
 120 measures the geometric similarity of representation spaces but also suffers from a high false positive  
 121 rate. More recently, TensorGuard (Wu et al., 2025) proposes utilizing the statistical features of gra-  
 122 dients generated during backpropagation as a stable signature to characterize the model’s optimization  
 123 landscape.

124 *Black-box Fingerprinting.* With the proliferation of Model-as-a-Service (MaaS), black-box meth-  
 125 ods that rely solely on API input-output interactions have gained traction. These approaches gen-  
 126 erally follow two paradigms: analyzing output distributions or constructing specific trigger queries.  
 127 The former focuses on identifying unique stylistic idiosyncrasies or probability distributions in-  
 128 herent to a model family. For example, LLMmap (Pasquini et al., 2025) and other distributional  
 129 approaches (Yang & Wu, 2024; Gao et al., 2025; McGovern et al., 2025) query the model with  
 130 general prompts to capture distinct response patterns or logit features. The latter paradigm involves  
 131 optimizing specific “trap” inputs or adversarial queries (Xu et al., 2024; Gubri et al., 2024; Jin et al.,  
 132 2024; Xu et al., 2025) designed to force a suspect model to output a predefined unique response.

133 However, black-box methods face significant robustness challenges. Fundamentally, relying solely  
 134 on surface-level outputs results in a loss of critical information regarding the model’s internal mech-  
 135 anisms compared to white-box access (Shao et al., 2025). Consequently, these methods are often  
 136 fragile to post-training modifications such as instruction tuning or simple system prompt changes,  
 137 which can disrupt the fingerprint (Xu et al., 2025; Tsai et al., 2025). To date, there is still a lack of  
 138 fingerprinting methods for LLMs that are both robust to extensive training and exhibit a very low  
 139 risk of false positives.

### 140 3 PRELIMINARY

141 **LLM Architecture** Most Large Language Models (LLMs) follows a decoder-only Transformer  
 142 architecture (Radford et al., 2018). The details of various LLMs may differ, yet the Transformer  
 143 blocks are similar. In particular, a Transformer block in an LLM usually consists of residual connec-  
 144 tions (He et al., 2016), Root Mean Square Normalization (RMSNorm, Zhang & Sennrich (2019)),  
 145 the self-attention mechanism (Lin et al., 2017), Rotary Position Embedding (RoPE, Su et al. (2024)),  
 146 and a feed-forward network (FFN). We denote by  $\mathbb{W}_A$  the set of an  $L$ -layered LLM A’s weights, and  
 147  $\mathbb{W}_{A,\text{partial}} = \{W_{A,\text{emb}}\} \bigcup_{l=1}^L \{W_{A,i}^{(l)} \mid i \in \{Q, K\}\}$  the set of word embeddings and Q, K matrices,  
 148 where  $W_{A,\text{emb}}$  is the word embeddings and  $W_{A,Q}^{(l)}, W_{A,K}^{(l)}$  are the Q, K matrices at the  $l$ -th layer. We  
 149 defer the rest of the notations to Appendix D.1.

150 **Central Kernel Alignment (CKA)** Proposed in Kornblith et al. (2019), CKA is a similarity detec-  
 151 tion method (Zhang et al., 2025) based on Hilbert-Schmidt Independence Criterion (HSIC, Gretton  
 152 et al. (2005)). It is invariant to column-wise orthogonal transformations and constant multiplications:

153 **Theorem 3.1** *For any input matrices  $X_1 \in \mathbb{R}^{m \times n_1}, X_2 \in \mathbb{R}^{m \times n_2}$ , any orthogonal transformations  
 154  $U_1 \in \mathbb{R}^{n_1 \times n_1}, U_2 \in \mathbb{R}^{n_2 \times n_2}$ , any non-zero constants  $c_1, c_2 \in \mathbb{R}$ ,*

$$155 \quad \text{CKA}(X_1, X_2) = \text{CKA}(c_1 X_1 U_1, c_2 X_2 U_2). \quad (1)$$

156 Even with semi-orthogonal  $U_1$  and  $U_2$  where  $U_1 \in \mathbb{R}^{n_1 \times n'_1}, U_2 \in \mathbb{R}^{n_2 \times n'_2}, n_1 > n'_1, n_2 > n'_2$ , CKA  
 157 still preserves input similarity to some extent (Kang et al., 2025; Chun et al., 2025). Nevertheless, the

standard HSIC estimator converges at rate  $1/\sqrt{m}$  and shows finite-sample bias (Gretton et al., 2005; Murphy et al., 2024). To this end, an unbiased version (Song et al., 2007) is proposed and extends the value of CKA from  $[0, 1]$  to  $[-1, 1]$ . Additionally, linear kernels are often selected in CKA due to their computational efficiency and similar performance to other kernels (Kornblith et al., 2019). We give the formal definition of CKA in Appendix D.2, and the proof of Theorem 3.1 in Appendix E.1.

## 4 MANIPULATIONS ON AN OPEN-SOURCE LLM

In this part, we investigate which matrix-weight manipulations can remain compatible with preserving a model’s behaviour. We first present the definition of matrix weight manipulation in §4.1. Then, we examine how key components of an LLM, including residual connections (§4.2), RMSNorm (§4.3), and RoPE together with attention scores (§4.4), can constrain and gradually narrow the space of admissible manipulations when the manipulated model is required to produce outputs similar to the base model. Finally, we derive potential attacks on matrix weights in §4.5 based on these constraints, with a focus on the Q, K matrices and the embedding layer.

### 4.1 PROBLEM DEFINITION

Weights of an open-source LLM are often inherited, but unclaimed inheritance invites manipulations. Although the source code may not explicitly reveal these manipulations, the weights are vulnerable to modifications that leave no trace in code (e.g., post-hoc matrix multiplications to produce new weights), and they can even be altered to evade independence tests. To enable detection on model weights, we first formalize the plausible manipulation forms.

**Definition 4.1 (LLM Weight Manipulations)** *Let  $A$  and  $B$  be two open-source LLMs with the same number of layers  $L$ , and  $\mathbb{W}_A, \mathbb{W}_B$  denote the sets of matrix weights for  $A$  and  $B$  respectively. Then, if  $B$  manipulates  $A$ , the manipulations on  $\mathbb{W}_{A,\text{partial}}$  are*

$$W_{B,i}^{(l)} = L_{B,i}^{(l)} W_{A,i}^{(l)} R_{B,i}^{(l)} + E_{B,i}^{(l)} \quad \text{for all } 1 \leq l \leq L \text{ and } i \in \{Q, K\} \quad (2)$$

$$W_{B,\text{emb}} = W_{A,\text{emb}} R_{B,\text{emb}} + E_{B,\text{emb}} \quad (3)$$

where  $L$  and  $R$  are row-wise and column-wise transformation matrices, and  $E$  is the error term. The learnable weights in RMSNorm are also changed correspondingly.

We omit the row-wise transform on word embeddings,  $L_{B,\text{emb}}$ , since each row of  $W_{A,\text{emb}}$  represents a token and mixed token representations are hard to be faithfully recovered in the calculation of attention scores. We further assume that the suspicious and manipulated models produce similar (or identical) outputs to be consistent with the goal of reusing base model performance. This assumption induces constraints on the transformation matrices and, in turn, enables detection via weight-matrix similarity tests. We first develop a fine-grained view of the constraints on  $\mathbb{W}_{A,\text{partial}}$  in what follows.

### 4.2 RESIDUAL CONNECTIONS: PASSING MANIPULATIONS FORWARD

In what follows, we first analyze how residual connections propagate weight manipulations forward through the network while preserving the model’s outputs. A Transformer block consists of multiple components linked by residual connections. If one component takes input  $X$ , the next receives  $Y = X + f(X)$ , so any manipulation  $T$  on  $X$  must be recovered within the component to propagate as the next input (Zeng et al., 2024). This propagation also depends on the constituent functions of the component, which either commute with the manipulation to recover it or remain invariant and absorb it. We formalize this in the following:

**Proposition 4.2 (Proof in Appendix E.2)** *Let  $f = f_n \circ \dots \circ f_1$  be a component in Transformer and let  $T$  be a manipulation in the input. If for every  $k \in \{1, \dots, n\}$ ,  $T \circ f_k = f_k \circ T$ , then  $f \circ T = T \circ f$ , i.e. the manipulation propagates through constituent functions of a component.*

Compared to Zeng et al. (2024), Proposition 4.2 shows that residual connections are more vulnerable under a component-wise view, since constituent functions can propagate manipulations in various ways. However, nonlinearities of the functions, especially those within the self-attention mechanism, pose constraints for manipulations on both inputs and weights. We next split the self-attention

216 mechanism into two constituent functions, RMSNorm and attention (see [Definition D.4](#)), and show  
 217 how these effects arise and are constrained these two functions with a focus on RMSNorm, RoPE  
 218 and attention scores.  
 219

#### 220 4.3 RMSNORM: A CONSTRAINT ON EMBEDDING MANIPULATIONS 221

222 [Next, we examine how RMSNorm may allow certain manipulations of word embeddings.](#) Any input  
 223 to the self-attention mechanism, including word embeddings and hidden states, first go through  
 224 RMSNorm. However, RMSNorm can facilitate potential manipulations on inputs, because it com-  
 225 mutes with certain transformations  $R_{B,\text{emb}}$  of the embeddings:  
 226

227 **Theorem 4.3** *Let models A and B share the same<sup>1</sup> architecture. Let  $c \neq 0$  be a scalar;  $P$  be a  
 228 (partial) permutation matrix, and  $D$  be a signature matrix. Then,  $R_{B,\text{emb}} = cPD$  can be recovered  
 229 after RMSNorm in model B if related RMSNorm parameters are adjusted.*  
 230

231 [Theorem 4.3](#) indicates that RMSNorm is susceptible to embedding manipulations composed of con-  
 232 stant multiplications, permutations and sign flips. On the other hand, other manipulations on word  
 233 embeddings are generally neither commutative with RMSNorm nor invariant to it, which potentially  
 234 brings a constraint to manipulations. We provide a proof for [Theorem 4.3](#) in [Appendix E.3](#), along  
 235 with a discussion on other manipulations on word embeddings.  
 236

#### 237 4.4 ROPE AND ATTENTION SCORES: BOUNDARIES FOR Q/K MANIPULATIONS 238

239 [Finally, we study RoPE and the attention score computation to characterize which manipulations  
 240 on Q,K matrices can preserve attention scores.](#) After RMSNorm, inputs are fed into attention score  
 241 calculations. The nonlinearity here, particularly in the softmax and RoPE functions, poses a barrier  
 242 for manipulations on inputs to pass through. Hence, we follow [Zeng et al. \(2024\)](#) to assume that input  
 243 manipulations does not change the value of attention scores. The manipulations on Q,K matrices  
 244 are thus constrained under this assumption.  
 245

246 **Theorem 4.4** *The manipulations on Q,K matrices at layer l can be categorized into*  
 247

- 248 *1. Input-related ones, passed by RMSNorm:  $W_{B,i}^{(l)} = c^{-1}W_{A,i}^{(l)}PD$ , for  $i \in \{Q,K\}$ ;*  
 249
- 250 *2. RoPE-related ones:  $W_{B,i}^{(l)} = U_{B,i}^{(l)}W_{A,i}^{(l)}$ , for  $i \in \{Q,K\}$ ;*  
 251

252 *where  $U_{B,i}^{(l)}$  are special orthogonal matrices that keep RoPE results.*  
 253

254 A proof for [Theorem 4.4](#) is provided in [Appendix E.4](#), where we also show how manipulations on  
 255 inputs are recovered after V, O matrices to satisfy [Proposition 4.2](#). These manipulations preserve the  
 256 attention score values of model A in the suspicious model B. However, they can greatly change the  
 257 weights of the original model A, bringing difficulty to the development of detection methods.  
 258

#### 259 4.5 POTENTIAL ATTACKS 260

261 Combining [Theorem 4.3](#) and [Theorem 4.4](#), the admissible manipulations on  $\mathbb{W}_{A,\text{partial}}$  (word embed-  
 262 dings and Q, K matrices) in [Definition 4.1](#) are restricted to  
 263

$$264 W_{B,i}^{(l)\top} = c^{-1} D^\top P^\top W_{A,i}^{(l)\top} U_{B,i}^{(l)\top} + E_{B,i}^{(l)\top}, \quad 1 \leq l \leq L, \quad i \in \{Q, K\}, \quad (4)$$

$$265 W_{B,\text{emb}} = c W_{A,\text{emb}} P D + E_{B,\text{emb}}. \quad (5)$$

266 Here  $E$  collects post-training, continual pre-training, pruning, upcycling, multimodal adaptation, or  
 267 related adjustments. The nonlinear usage of  $\mathbb{W}_{A,\text{partial}}$ , especially through the Q, K matrices, sub-  
 268 stantially limits manipulation complexity. Consequently, checking similarity between  $\mathbb{W}_{A,\text{partial}}$  and  
 269  $\mathbb{W}_{B,\text{partial}}$  is typically sufficient for detection. The remaining avenues targeting  $\mathbb{W}_A/\mathbb{W}_{A,\text{partial}}$  are de-  
 270 ferred to [Appendix E.5](#), where we also discuss how such transformations can recover Transformer  
 271 block outputs in light of [Proposition 4.2](#).  
 272

273 <sup>1</sup>The conclusions generalize to models with different number of layers (possibly a result of pruning).  
 274 See [Appendix E.3](#)

270 5 METHODOLOGY  
271272 Building on [Theorem 4.3](#) and [Theorem 4.4](#), we propose a two-stage procedure in [Algorithm 1](#): (i)  
273 extract  $P$  and  $D$  from the word embeddings; (ii) assess cross-model similarity by comparing the  $Q$   
274 and  $K$  matrices. This design achieves fast, reliable discrimination with minimal computation.  
275276 **Algorithm 1** LAP-Enhanced Unbiased Central Kernel Alignment (UCKA) Similarity Detection  
277

```

278 1: Given two  $L^2$ -layered LLMs  $A, B$  and their weight matrices  $\mathbb{W}_{A,\text{partial}}, \mathbb{W}_{B,\text{partial}}$ .
279 2: Let  $I = \text{Vocab}(A) \cap \text{Vocab}(B)$ , and  $m' = |I|$ .
280 3: Let  $W_{A,\text{shared-emb}} = W_{A,\text{emb}}[I, :]$  and  $W_{B,\text{shared-emb}} = W_{B,\text{emb}}[I, :]$ .
281 4: Build the cosine similarity matrix  $C$  with  $C_{k,l} = \frac{\langle (W_{A,\text{shared-emb}})_{:,k}, (W_{B,\text{shared-emb}})_{:,l} \rangle}{\|(W_{A,\text{shared-emb}})_{:,k}\| \|(W_{B,\text{shared-emb}})_{:,l}\|}$ .
282 5: Construct the permutation and signature matrices  $P, D$  via LAP:
283 6:   Find a permutation  $\pi$  maximizing  $\sum_k |C_{k,\pi(k)}|$  with the Hungarian algorithm.
284 7:   For each column  $k$ , set  $s_k = \text{sign}(C_{k,\pi(k)})$ .
285 8:   Set  $P_{k,\pi(k)} = 1$  for every  $k$ , set  $D = \text{diag}(s_1, s_2, \dots, s_k, \dots)$ .
286 9: for  $l = 1, \dots, L$  do
287 10:    $s_Q^{(l)} = \text{UCKA}\left(D^\top P^\top W_{A,Q}^{(l)\top}, W_{B,Q}^{(l)\top}\right)$ ,  $s_K^{(l)} = \text{UCKA}\left(D^\top P^\top W_{A,K}^{(l)\top}, W_{B,K}^{(l)\top}\right)$ 
288 11: end for
289 12: return  $\sum_{l=1}^L (|s_Q^{(l)}| + |s_K^{(l)}|)/2L$ 
290
291
292
```

293 **Extracting the Permutation and Signature from Word Embeddings** A key property of  $R_{B,\text{emb}}$   
294 is that the permutation and signature matrices may appear in either order: both  $R_{B,\text{emb}} = cPD$   
295 and  $R_{B,\text{emb}} = cDP$  are possible. However, any product  $DP$  can be rewritten as  $P'D'$  for some  
296 permutations  $P'$  and signature matrices  $D'$ . Hence it suffices to recover a canonical  $PD$  from the  
297 embeddings. We cast this as a Linear Assignment Problem (LAP; [Burkard & Cela \(1999\)](#)) solved by  
298 the Hungarian algorithm ([Kuhn, 1955](#)), which is invariant to any nonzero scalar factor. Specifically,  
299 we first restrict to the shared vocabulary and form the matrix of absolute cosine similarities between  
300 the columns of  $W_{A,\text{emb}}$  and  $W_{B,\text{emb}}$ . Then, LAP is applied to obtain the permutation  $P$ . Next, we use  
301 the signs of the cosine similarities at the matched pairs to reconstruct the signature matrix  $D$ . This  
302 approach effectively reconstruct corresponding transformations due to its low demand of additional  
303 parameters in the detection process.  
304305 **Robust Recovery of Weight Similarities** Despite accounting for permutation and signature ma-  
306 nipulations, the orthogonal transformations  $U_{B,i}^{(l)}$  remain challenging. They introduce a substantial  
307 nuisance parameter burden: if  $W_{A,i}^{(l)} \in \mathbb{R}^{d \times n}$ , then orthogonal  $U_{B,i}^{(l)} \in \mathbb{R}^{d \times d}$  contributes  $d^2$  pa-  
308 rameters, which is prohibitive when  $d \approx n$  and can undermine robustness by adding parameters to de-  
309 tections ([Simmons et al., 2011](#)). We therefore use Central Kernel Alignment (CKA) as a parameter-free  
310 similarity metric: by [Theorem 3.1](#), CKA is invariant to orthogonal transforms and constant rescal-  
311 ing. While this invariance does not extend to semi-orthogonal  $U_{B,i}^{(l)}$  (e.g., induced by pruning), CKA  
312 remains effective to a meaningful extent in such cases, relieving the need to explicitly reconstruct  
313  $U_{B,i}^{(l)}$ . To further mitigate biases ([Gretton et al., 2005](#); [Murphy et al., 2024](#)), we adopt the unbiased  
314 variant ([Song et al., 2007](#)), termed UCKA (see [Appendix D.2](#)).316 6 EXPERIMENTS  
317318 In this section, we present a series of experiments designed to rigorously evaluate our proposed  
319 model fingerprinting method. We begin in [Section 6.1](#) by verifying its fundamental ability to dis-  
320 tinguish between derived (offspring) and independent models. Next, in [Section 6.2](#), we assess the  
321 critical risk of false positives by comparing its performance on 90 pairs of independent models  
322323 <sup>2</sup>For LLMs with differing layer counts ( $L_A, L_B$ ), e.g., from layer pruning, we find the optimal layer pairing  
by solving the LAP on an  $L_A \times L_B$  matrix of layer-wise similarities before calculating overall similarity.

324  
 325  
 326  
 327  
 328 Table 1: Similarity (%) of various LLMs to LLaMA2-7B and LLaMA2-13B base models. Offspring  
 329 models consistently show high similarity scores (light red), while independent models have negligi-  
 330 ble similarity (light green), demonstrating the method’s discriminative power.  
 331  
 332  
 333  
 334  
 335  
 336  
 337

| Base Model: LLaMA2-7B |       |              |      | Base Model: LLaMA2-13B |       |               |      |
|-----------------------|-------|--------------|------|------------------------|-------|---------------|------|
| Offspring             | Sim   | Independent  | Sim  | Offspring              | Sim   | Independent   | Sim  |
| WizardMath-7b         | 99.99 | Baichuan-7b  | 0.58 | Selfrag_llama2_13b     | 99.99 | Baichuan-13b  | 0.17 |
| Selffrag_7b           | 99.98 | Mistral-7b   | 0.63 | Nous-Hermes-13b        | 99.98 | Baichuan2-13b | 0.23 |
| Vicuna-7b             | 99.95 | OLMo-7b      | 0.46 | Llama2-13b-orca        | 99.98 | OLMo2-13b     | 0.21 |
| Llama2-7b-Chat        | 99.93 | Qwen-7b      | 0.19 | Vicuna-13b             | 99.93 | Qwen-14b      | 0.20 |
| Finance-7b            | 99.93 | InternLM-7b  | 0.43 | Llama2-13b-Chat        | 99.92 | Qwen3-14b     | 0.41 |
| Llama2-7b-32K         | 99.88 | MPT-7b       | 0.04 | Firefly-llama2-13b     | 99.81 | Jais-13b      | 0.03 |
| Guanaco-7b            | 99.80 | LLaMA-7b     | 0.81 | Llama2-koen-13b        | 97.76 | LLaMA-13b     | 0.74 |
| Llama2-ko-7b          | 96.79 | OpenLlama-7b | 0.71 | Llama2-13b-Estopia     | 96.60 | OpenLlama-13b | 0.51 |

338  
 339 against two state-of-the-art baselines, HuRef and REEF. In [Section 6.3](#), we test the method’s robust-  
 340 ness against a wide array of common post-training modifications using a comprehensive suite of 60  
 341 offspring models. Finally, in [Section 6.4](#), we provide an overall performance comparison, leveraging  
 342 ROC curves and other metrics to demonstrate our method’s superior discriminative power.  
 343

344 **Baselines.** We compare our method against two advanced baselines: the weight-based method  
 345 HuRef ([Zeng et al., 2024](#)), and the representation-based method REEF ([Zhang et al., 2025](#)).  
 346

## 347 6.1 EFFECTIVENESS VERIFICATION

348 We first established our method’s core effectiveness in identifying model lineage. We collected eight  
 349 offspring models for both LLaMA2-7B and LLaMA2-13B, alongside eight independent models for  
 350 each size. We then calculated the similarity between the base model and these two groups.  
 351

352 As shown in [Table 1](#), the results demonstrate a clear distinction. Offspring models exhibited ex-  
 353 ceptionally high similarity to their respective base models, whereas the similarity scores between  
 354 independent models were negligible, forming a strong basis for intellectual property protection.  
 355

## 356 6.2 FALSE POSITIVE RISK EVALUATION ON 90 UNRELATED PAIRS

357 A critical requirement for any reliable fingerprinting method is an extremely low false positive rate.  
 358 To evaluate this risk, we collected 10 independent 7B LLMs and 10 independent 13B LLMs, forming  
 359 45 unique pairs for each size. We then computed the pairwise similarity (%) for all 90 pairs using  
 360 our method and compared the results with those from REEF and HuRef.  
 361

362 The heatmaps in [Figure 2](#) illustrate our method’s superior performance in avoiding false positives.  
 363 For independent pairs, our method yielded mean similarity scores of just 0.49 (7B) and 0.26 (13B).  
 364 These scores are nearly an order of magnitude lower than HuRef (means of 3.56 and 2.17) and two  
 365 orders of magnitude lower than REEF (means of 42.47 and 47.44).  
 366

367 Notably, REEF frequently produces dangerously high similarity scores for unrelated models, with  
 368 many pairs exceeding 80 and some even surpassing 95. Such high values could easily lead to false  
 369 accusations of model theft. In contrast, the maximum similarity scores for our method were merely  
 370 1.5 (7B) and 0.8 (13B), reaffirming its significantly lower risk of false positives.  
 371

## 372 6.3 ROBUSTNESS VERIFICATION ON 60 OFFSPRING MODELS

373 Base LLMs often undergo substantial post-training modifications, including SFT, continued pre-  
 374 training, reinforcement learning (RL), pruning, upcycling, and multi-modal adaptation. These pro-  
 375 cesses can significantly alter model parameters, making robustness a critical attribute for any finger-  
 376 printing technique. To rigorously evaluate our method’s robustness, we curated a diverse test suite of  
 377 60 positive pairs, each consisting of a base model and a derived offspring model. This suite includes  
 378 10 pairs for each of the six modification categories listed above.  
 379

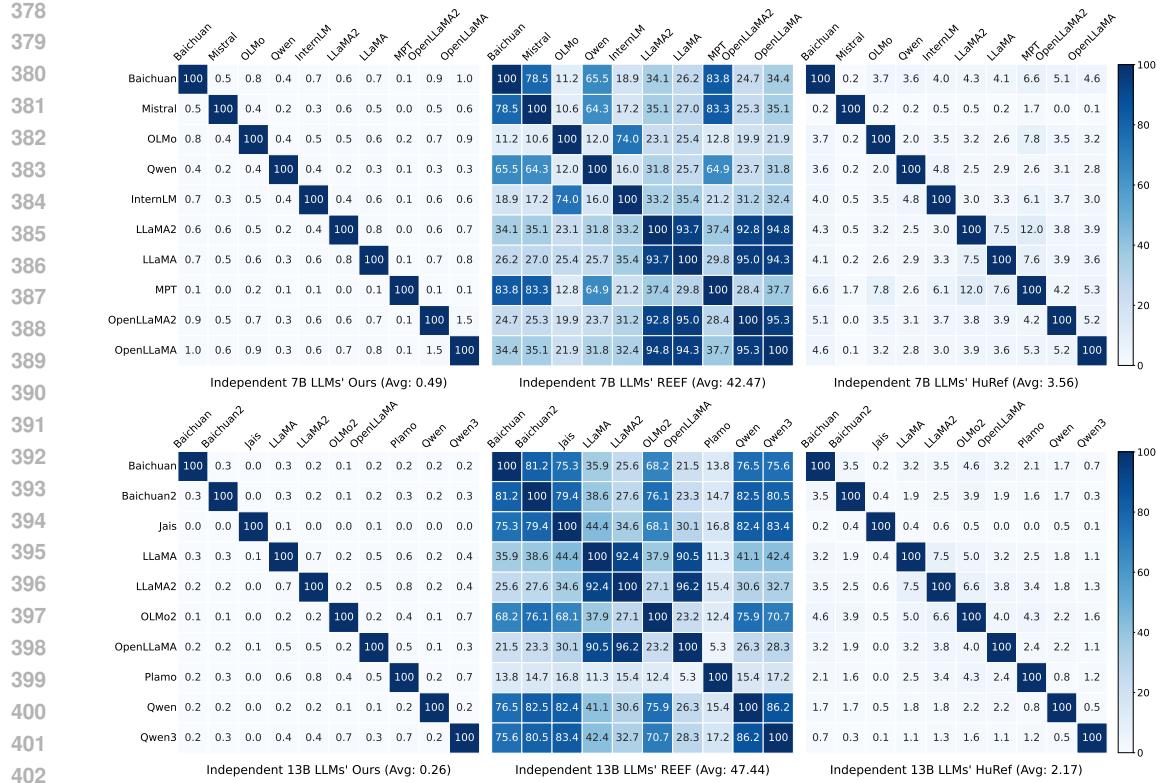


Figure 2: Pairwise similarity (%) heatmaps for independent 7B (top row) and 13B (bottom row) LLMs. The columns compare our proposed method against two baselines, REEF and HuRef. Our method consistently yields near-zero similarity scores for all independent model pairs, indicating a significantly lower risk of false positives. In contrast, REEF often produces high similarity scores ( $> 80$ ), which could easily lead to false accusations of model theft.

The effectiveness of a fingerprinting method hinges on its ability to distinguish between related (positive) and unrelated (negative) pairs. However, absolute similarity scores (%) can be misleading. A method yielding scores of 10 for a positive pair and 0.1 for negative pairs is far more discriminative than one producing scores of 90 and 40. To accurately measure the discriminative power, we must quantify how statistically different a positive pair’s similarity is from the distribution of negative pairs. For this, we use the absolute Z-score ( $|Z|$ )<sup>3</sup>. This metric measures how many standard deviations a positive pair’s similarity is from the mean of the negative pairs’ similarities. A larger  $|Z|$  indicates that the positive pair is more likely to be a statistical outlier relative to the population of unrelated pairs and thus more highly separable.

The results, presented in Table 2, highlight the superior robustness of our method. HuRef loses effectiveness under heavy modifications like extensive continued pretraining (e.g., its  $|Z|$  drops to 0.57 for Qwen2.5-coder). REEF consistently yields low  $|Z|$ , often below 2.0, indicating poor separability. In stark contrast, our method maintains remarkably high  $|Z|$  across all 60 positive pairs, demonstrating its resilience against a wide array of demanding post-training modifications.

#### 6.4 OVERALL EVALUATION

To synthesize our findings, we conducted a comprehensive evaluation against HuRef (ICS) and REEF using the full dataset of 60 positive and 90 negative model pairs. Additionally, we incorporated PCS (from HuRef) and Intrinsic Fingerprint (Yoon et al., 2025) as supplementary baselines to broaden the comparative analysis.

<sup>3</sup>In this context, the absolute Z-score is mathematically equivalent to the Mahalanobis distance.

432 Table 2: Absolute Z-scores ( $\uparrow$ ) of positive samples for HuRef, REEF, and our method. The full  
 433 names of abbreviations and related model details are provided in [Appendix G](#).  
 434

| SFT                |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
|--------------------|--------------------|--------------------|------------------|----------------|-------------------|-------------------|-----------------|-------------------|-----------------|-------------------|
| Base               | Llama2-7B          |                    |                  |                |                   | Llama2-13B        |                 |                   |                 |                   |
| Offspring          | Vicuna             | Finance            | Selfrag          | 32K            | Wizard            | Guanaco           | Vicuna          | Hermes            | Estopia         | Firefly           |
| HuRef              | 44.22              | 44.20              | 44.47            | 44.22          | 44.51             | 44.12             | 44.25           | 44.48             | 43.09           | 44.51             |
| REEF               | 1.95               | 1.94               | 1.95             | 1.88           | 1.95              | 1.94              | 1.94            | 1.95              | 1.91            | 1.95              |
| Ours               | 355.09             | 355.02             | 355.20           | 354.84         | 355.23            | 354.55            | 355.02          | 355.20            | 343.14          | 354.59            |
| Continual Pretrain |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
| Base               | Llama2-7B          |                    |                  | Gemma          | Gemma             | Qwen2.5           |                 | Qwen2             | Llama2-70B      |                   |
| Offspring          | Llemma<br>(700B)   | Code<br>(520B)     | Python<br>(620B) | Code<br>(500B) | Code<br>(500B)    | Math<br>(1T)      | Coder<br>(5.5T) | Math<br>(700B)    | Code<br>(520B)  | Python<br>(620B)  |
| HuRef              | 8.15               | 9.41               | 9.19             | 28.97          | 39.51             | 0.28              | 0.57            | 1.01              | 3.43            | 2.78              |
| REEF               | 1.66               | 1.31               | 1.77             | 0.32           | 1.52              | 1.04              | 1.29            | 1.09              | 1.93            | 1.92              |
| Ours               | 250.32             | 250.78             | 253.28           | 198.18         | 268.72            | 177.57            | 135.17          | 183.49            | 241.12          | 232.92            |
| Upcycling          |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
| Base               | Mistral<br>-7B     | Llama3<br>-8B      |                  | Llama2-7B      |                   |                   |                 |                   | Qwen<br>-1.8B   | Minicpm<br>-2B    |
| Offspring          | Mixtral            | MoE<br>v2          | MoE4             | MoE<br>3B      | MoE2              | MoE3B<br>-SFT     | MoE2<br>-SFT    | MoE4<br>-SFT      | Qwen1.5<br>MoE  | Minicpm<br>MoE    |
| HuRef              | 7.97               | 42.50              | 24.58            | 21.59          | 24.91             | 21.59             | 24.91           | 24.58             | 4.03            | 11.39             |
| REEF               | 1.47               | 0.47               | 0.64             | 0.64           | 0.62              | 0.63              | 0.63            | 0.62              | 0.57            | 1.49              |
| Ours               | 239.01             | 332.12             | 332.23           | 326.49         | 332.12            | 326.49            | 332.12          | 332.23            | 173.36          | 150.15            |
| Multi Modal        |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
| Base               | Llama2-7B          |                    | Qwen2<br>-7B     | Qwen-7B        |                   |                   | Qwen2.5<br>-7B  | Qwen2.5<br>-3B    | Llama3<br>-8B   | Llama2<br>-13B    |
| Offspring          | LLaVA              | Video              | VL               | Audio          | Audio2            | VL                | VL              | VL                | Next            | LLaVA             |
| HuRef              | 44.06              | 44.03              | 39.94            | 38.79          | 21.59             | 37.41             | 24.28           | 23.95             | 44.30           | 44.09             |
| REEF               | 0.40               | 0.36               | 0.84             | 0.19           | 0.45              | 0.77              | 0.48            | 0.61              | 0.26            | 0.07              |
| Ours               | 354.98             | 354.98             | 342.72           | 339.86         | 317.97            | 336.55            | 290.08          | 298.46            | 355.05          | 354.91            |
| RL                 |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
| Base               | Open-<br>llama3B   | Qwen2.5<br>-7B     | Qwen2.5<br>-1.5B | Mixtral        | Mistral-7B        |                   | Minicpm<br>-2B  | Qwen3<br>-4B      | Chatglm<br>-6B  | Llama3<br>-8B     |
| Offspring          | RLHF               | Reason             | Zero             | DPO            | DPO               | Dolphin           | DPO             | GRPO              | RLHF            | DPO               |
| HuRef              | 44.52              | 44.58              | 44.58            | 44.57          | 44.53             | 44.54             | 44.58           | 44.58             | 44.58           | 44.52             |
| REEF               | 1.94               | 1.93               | 1.94             | 1.75           | 1.48              | 1.21              | 1.92            | 1.78              | 1.96            | 1.96              |
| Ours               | 355.23             | 355.27             | 355.27           | 355.23         | 355.23            | 355.23            | 355.27          | 355.27            | 355.27          | 355.23            |
| Pruning            |                    |                    |                  |                |                   |                   |                 |                   |                 |                   |
| Base               | Llama-3-8B         |                    |                  |                |                   | Llama2-7B         |                 |                   |                 |                   |
| Offspring          | Minitron<br>-Depth | Minitron<br>-Width | Llama3<br>-1B    | Llama3<br>-3B  | Sheared<br>2.7B-P | Sheared<br>2.7B-S | Sheared<br>2.7B | Sheared<br>1.3B-P | Sheared<br>1.3B | Sheared<br>1.3B-S |
| HuRef              | 28.29              | 22.23              | 0.33             | 0.73           | 22.88             | 16.00             | 15.86           | 10.06             | 7.64            | 7.64              |
| REEF               | 0.52               | 0.53               | 1.15             | 0.92           | 1.75              | 1.78              | 1.77            | 1.80              | 1.79            | 1.79              |
| Ours               | 344.07             | 343.00             | 12.14            | 106.29         | 328.81            | 312.44            | 312.80          | 317.79            | 297.50          | 297.50            |

485 As illustrated in [Figure 1](#), our method demonstrates vastly superior performance. The Receiver  
 486 Operating Characteristic (ROC) curve (left) shows that our method achieves a perfect Area Under

Table 3: Detailed performance comparison of fingerprinting methods across various post-training techniques. Our method consistently achieves perfect scores (1.0) on all classification metrics (AUC, pAUC, TPR@1%FPR) and maintains a significantly larger separation margin ( $|\bar{Z}|$ ) across all scenarios. CPT: Continual Pre-Training, UP: Upcycling, MM: Multi-modal, PR: Pruning.

| Method                | Metric               | SFT     | CPT     | UP      | MM      | RL      | PR      | All     |
|-----------------------|----------------------|---------|---------|---------|---------|---------|---------|---------|
| HuRef                 | $ \bar{Z}  \uparrow$ | 43.748  | 10.331  | 20.805  | 36.244  | 44.559  | 13.166  | 28.142  |
|                       | AUC $\uparrow$       | 1.000   | 0.879   | 0.999   | 1.000   | 1.000   | 0.866   | 0.957   |
|                       | pAUC $\uparrow$      | 1.000   | 0.656   | 0.978   | 1.000   | 1.000   | 0.800   | 0.906   |
|                       | TPR@1%FPR $\uparrow$ | 1.000   | 0.500   | 0.900   | 1.000   | 1.000   | 0.800   | 0.867   |
| REEF                  | $ \bar{Z}  \uparrow$ | 1.936   | 1.384   | 0.777   | 0.443   | 1.788   | 1.381   | 1.285   |
|                       | AUC $\uparrow$       | 1.000   | 0.857   | 0.508   | 0.648   | 0.963   | 0.692   | 0.778   |
|                       | pAUC $\uparrow$      | 1.000   | 0.211   | 0.000   | 0.000   | 0.658   | 0.300   | 0.362   |
|                       | TPR@1%FPR $\uparrow$ | 1.000   | 0.200   | 0.000   | 0.000   | 0.600   | 0.000   | 0.300   |
| Intrinsic Fingerprint | $ \bar{Z}  \uparrow$ | 1.535   | 1.193   | 1.408   | 1.532   | 1.542   | 1.141   | 1.392   |
|                       | AUC $\uparrow$       | 1.000   | 0.896   | 0.969   | 1.000   | 1.000   | 0.876   | 0.957   |
|                       | pAUC $\uparrow$      | 1.000   | 0.422   | 0.800   | 1.000   | 1.000   | 0.400   | 0.770   |
|                       | TPR@1%FPR $\uparrow$ | 1.000   | 0.300   | 0.800   | 1.000   | 1.000   | 0.400   | 0.750   |
| PCS                   | $ \bar{Z}  \uparrow$ | 73.786  | 74.650  | 0.727   | 14.950  | 163.399 | 17.226  | 57.457  |
|                       | AUC $\uparrow$       | 0.958   | 0.959   | 0.791   | 0.802   | 0.984   | 0.666   | 0.860   |
|                       | pAUC $\uparrow$      | 0.656   | 0.600   | 0.078   | 0.500   | 0.900   | 0.100   | 0.472   |
|                       | TPR@1%FPR $\uparrow$ | 0.500   | 0.600   | 0.000   | 0.500   | 0.900   | 0.100   | 0.433   |
| Ours                  | $ \bar{Z}  \uparrow$ | 353.788 | 219.155 | 287.634 | 334.556 | 355.250 | 267.233 | 302.936 |
|                       | AUC $\uparrow$       | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   |
|                       | pAUC $\uparrow$      | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   |
|                       | TPR@1%FPR $\uparrow$ | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   | 1.000   |

the Curve (AUC) of 1.0, indicating flawless discrimination. This starkly outperforms both HuRef ( $AUC = 0.957$ ) and REEF ( $AUC = 0.778$ ).

In practical applications, preventing false accusations of model theft is paramount, making performance at a very low False Positive Rate (FPR) essential. We therefore employ two stricter metrics: the partial AUC for  $FPR < 0.05$  (pAUC) and the True Positive Rate at a 1% FPR (TPR@1%FPR).

**Table 3** details the performance breakdown across post-training techniques. Our method consistently achieves perfect scores (1.0) on all classification metrics (AUC, pAUC, and TPR@1%FPR) across all categories. In contrast, the baseline methods show significant limitations. REEF’s performance collapses in several scenarios, with pAUC and TPR@1%FPR scores falling to 0.0 for Upcycling and Multi-modal. HuRef also shows vulnerability, with its TPR@1%FPR dropping to 0.500 under Continual Pre-Training. Our method’s perfect classification scores, combined with its substantially larger average absolute Z-score ( $|\bar{Z}|$ ), underscore its superior robustness and reliability.

## 7 CONCLUSIONS

In this paper, we propose a training-free fingerprinting method for LLM identification. Our approach does not impair LLM’s general capability while exhibiting robustness against fine-tuning, extensive continued pretraining, reinforcement learning, multimodal extension, pruning, and upcycling, and simultaneously avoids the risk of false positives. Experiments on a testbed comprising 150 pairs of LLMs demonstrate the effectiveness of our method.

540 REFERENCES  
541

542 Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. Turning your weak-  
543 ness into a strength: Watermarking deep neural networks by backdooring. In *27th USENIX Secu-  
544 rity Symposium (USENIX Security 18)*, pp. 1615–1631, 2018.

545 BaiChuan-Inc. <https://github.com/baichuan-inc/Baichuan-7B>, 2023.  
546

547 Rainer E Burkard and Eranda Cela. Linear assignment problems and extensions. In *Handbook of  
548 combinatorial optimization: Supplement volume A*, pp. 75–149. Springer, 1999.

549 Huili Chen, Bita Darvish Rouhani, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar. Deepmarks:  
550 A secure fingerprinting framework for digital rights management of deep learning models. In  
551 *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pp. 105–113,  
552 2019a.

553 Huili Chen, Bita Darvish Rouhani, and Farinaz Koushanfar. Blackmarks: Blackbox multibit water-  
554 marking for deep neural networks. *arXiv preprint arXiv:1904.00344*, 2019b.  
555

556 Jialuo Chen, Jingyi Wang, Tinglan Peng, Youcheng Sun, Peng Cheng, Shouling Ji, Xingjun Ma,  
557 Bo Li, and Dawn Song. Copy, right? a testing framework for copyright protection of deep  
558 learning models. In *IEEE Symposium on Security and Privacy*, pp. 824–841, 2022.  
559

560 Chanwoo Chun, Abdulkadir Canatar, SueYeon Chung, and Daniel D Lee. Estimating neural repre-  
561 sentation alignment from sparsely sampled inputs and features. *arXiv preprint arXiv:2502.15104*,  
562 2025.

563 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
564 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of  
565 the North American chapter of the association for computational linguistics: human language  
566 technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.  
567

568 Pierre Fernandez, Guillaume Couairon, Teddy Furon, and Matthijs Douze. Functional invariants to  
569 watermark large transformers. In *ICASSP 2024-2024 IEEE International Conference on Acous-  
570 tics, Speech and Signal Processing (ICASSP)*, pp. 4815–4819. IEEE, 2024.  
571

572 Irena Gao, Percy Liang, and Carlos Guestrin. Model equality testing: Which model is this api  
573 serving? In *International Conference on Learning Representations*, 2025.  
574

575 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
576 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd  
577 of models. *arXiv preprint arXiv:2407.21783*, 2024.  
578

579 Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical de-  
580 pendence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*,  
581 pp. 63–77. Springer, 2005.  
582

583 Martin Gubri, Dennis Ulmer, Hwaran Lee, Sangdoo Yun, and Seong Joon Oh. Trap: Targeted  
584 random adversarial prompt honeypot for black-box identification. In *Findings of the Association  
585 for Computational Linguistics: ACL 2024*, pp. 11496–11517, 2024.  
586

587 Jia Guo and Miodrag Potkonjak. Watermarking deep neural networks for embedded systems. In  
588 *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pp. 1–8. IEEE,  
589 2018.  
590

591 Ethan He, Abhinav Khattar, Ryan Prenger, Vijay Korthikanti, Zijie Yan, Tong Liu, Shiqing Fan,  
592 Ashwath Aithal, Mohammad Shoeybi, and Bryan Catanzaro. Upcycling large language models  
593 into mixture of experts. *arXiv preprint arXiv:2410.07524*, 2024.  
594

595 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
596 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
597 770–778, 2016.  
598

594 Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,  
 595 Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*,  
 596 2024.

597 Heng Jin, Chaoyu Zhang, Shanghao Shi, Wenjing Lou, and Y Thomas Hou. Proflingo: A  
 598 fingerprinting-based intellectual property protection scheme for large language models. In *IEEE*  
 599 *Conference on Communications and Network Security*, pp. 1–9. IEEE, 2024.

600 Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,  
 601 Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv*  
 602 *preprint arXiv:2503.19786*, 2025.

603 Hyunmo Kang, Abdulkadir Canatar, and SueYeon Chung. Spectral analysis of representational  
 604 similarity with limited neurons. *arXiv preprint arXiv:2502.19648*, 2025.

605 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A  
 606 watermark for large language models. In *International Conference on Machine Learning*, pp.  
 607 17061–17084. PMLR, 2023.

608 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural  
 609 network representations revisited. In *International conference on machine learning*, pp. 3519–  
 610 3529. PMIR, 2019.

611 Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics*  
 612 *quarterly*, 2(1-2):83–97, 1955.

613 Erwan Le Merrer, Patrick Perez, and Gilles Trédan. Adversarial frontier stitching for remote neural  
 614 network watermarking. *Neural Computing and Applications (NCA)*, 32(13):9233–9244, 2020.

615 Peixuan Li, Pengzhou Cheng, Fangqi Li, Wei Du, Haodong Zhao, and Gongshen Liu. Plmmmark:  
 616 a secure and robust black-box watermarking framework for pre-trained language models. In  
 617 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 14991–14999,  
 618 2023.

619 Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou,  
 620 and Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv preprint*  
 621 *arXiv:1703.03130*, 2017.

622 Hanwen Liu, Zhenyu Weng, and Yuesheng Zhu. Watermarking deep neural networks with greedy  
 623 residuals. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp.  
 624 6978–6988. PMLR, 2021.

625 Nils Lukas, Yuxuan Zhang, and Florian Kerschbaum. Deep neural network fingerprinting by con-  
 626 ferrable adversarial examples. *arXiv preprint arXiv:1912.00888*, 2019.

627 Hope Elizabeth McGovern, Rickard Stureborg, Yoshi Suhara, and Dimitris Alikaniotis. Your large  
 628 language models are leaving fingerprints. In *International Conference on Computational Linguis-  
 629 tics Workshop*, pp. 85–95, 2025.

630 Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent  
 631 Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based  
 632 on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.

633 Alex Murphy, Joel Zylberberg, and Alona Fyshe. Correcting biased centered kernel alignment  
 634 measures in biological and artificial neural networks. *arXiv preprint arXiv:2405.01012*, 2024.

635 Xudong Pan, Yifan Yan, Mi Zhang, and Min Yang. Metav: A meta-verifier approach to task-agnostic  
 636 model fingerprinting. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Dis-  
 637 covery and Data Mining (SIGKDD)*, pp. 1327–1336, 2022.

638 Dario Pasquini, Evgenios M Kornaropoulos, and Giuseppe Ateniese. Llmmmap: Fingerprinting for  
 639 large language models. In *USENIX Security Symposium*, 2025.

648 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,  
 649 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The RefinedWeb  
 650 dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv*  
 651 preprint [arXiv:2306.01116](https://arxiv.org/abs/2306.01116), 2023. URL <https://arxiv.org/abs/2306.01116>.

652 Wenjun Peng, Jingwei Yi, Fangzhao Wu, Shangxi Wu, Bin Bin Zhu, Lingjuan Lyu, Binxing Jiao,  
 653 Tong Xu, Guangzhong Sun, and Xing Xie. Are you copying my model? protecting the copyright  
 654 of large language models for easas via backdoor watermark. In *Proceedings of the 61st Annual*  
 655 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7653–  
 656 7668, 2023.

657 Zirui Peng, Shaofeng Li, Guoxing Chen, Cheng Zhang, Haojin Zhu, and Minhui Xue. Fingerprinting-  
 658 ing deep neural networks globally via universal adversarial perturbations. In *Proceedings of the*  
 659 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13430–13439,  
 660 2022.

661 pzc163 et al. Project author team stay tuned: I found out that the llama3-v project is stealing a lot  
 662 of academic work from minicpm-llama3-v 2.5. github issue opened in the minicpm-o repository.  
 663 <https://github.com/OpenBMB/MiniCPM-o/issues/196>, 2024. OpenBMB, June  
 664 2, 2024.

665 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language under-  
 666 standing by generative pre-training. 2018.

667 Mark Russinovich and Ahmed Salem. Hey, that's my model! introducing chain & hash, an llm  
 668 fingerprinting technique. *arXiv preprint arXiv:2407.10887*, 2024.

669 Shuo Shao, Yiming Li, Hongwei Yao, Yifei Chen, Yuchen Yang, and Zhan Qin. Reading between the  
 670 lines: Towards reliable black-box llm fingerprinting via zeroth-order gradient estimation. *arXiv*  
 671 preprint [arXiv:2510.06605](https://arxiv.org/abs/2510.06605), 2025.

672 Joseph P Simmons, Leif D Nelson, and Uri Simonsohn. False-positive psychology: Undisclosed  
 673 flexibility in data collection and analysis allows presenting anything as significant. *Psychological*  
 674 *science*, 22(11):1359–1366, 2011.

675 Le Song, Alex Smola, Arthur Gretton, Karsten M Borgwardt, and Justin Bedo. Supervised feature  
 676 selection via dependence estimation. In *Proceedings of the 24th international conference on*  
 677 *Machine learning*, pp. 823–830, 2007.

678 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-  
 679 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

680 CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu,  
 681 Christopher A Choquette-Choo, Jingyue Shen, Joe Kelley, et al. Codegemma: Open code models  
 682 based on gemma. *arXiv preprint arXiv:2406.11409*, 2024.

683 InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities.  
 684 <https://github.com/InternLM/InternLM>, 2023.

685 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen,  
 686 Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv*  
 687 preprint [arXiv:2507.20534](https://arxiv.org/abs/2507.20534), 2025.

688 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
 689 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
 690 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

691 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
 692 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
 693 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

694 Yun-Yun Tsai, Chuan Guo, Junfeng Yang, and Laurens van der Maaten. Rofl: Robust fingerprinting  
 695 of language models. *arXiv preprint arXiv:2505.12682*, 2025.

702 Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin'ichi Satoh. Embedding watermarks  
 703 into deep neural networks. In *Proceedings of the 2017 ACM on International Conference on*  
 704 *Multimedia Retrieval*. ACM, jun 2017. doi: 10.1145/3078971.3078974. URL <https://doi.org/10.1145%2F3078971.3078974>.

705

706 Tianhao Wang and Florian Kerschbaum. Riga: Covert and robust white-box watermarking of deep  
 707 neural networks. In *Proceedings of the Web Conference 2021 (WWW)*, pp. 993–1004, 2021.

708

709 Zehao Wu, Yanjie Zhao, and Haoyu Wang. Gradient-based model fingerprinting for llm similarity  
 710 detection and family classification. *arXiv preprint arXiv:2506.01631*, 2025.

711

712 Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language  
 713 model pre-training via structured pruning. In *The Twelfth International Conference on Learning*  
 714 *Representations*.

715

716 Jiašu Xu, Fei Wang, Mingyu Ma, Pang Wei Koh, Chaowei Xiao, and Muhaoo Chen. Instructional  
 717 fingerprinting of large language models. In *Proceedings of the 2024 Conference of the North*  
 718 *American Chapter of the Association for Computational Linguistics: Human Language Tech-*  
 719 *nologies (Volume 1: Long Papers)*, pp. 3277–3306, 2024.

720

721 Zhenhua Xu, Zhebo Wang, Maike Li, Wenpeng Xing, Chunqiang Hu, Chen Zhi, and Meng Han.  
 722 Rap-sm: Robust adversarial prompt via shadow models for copyright verification of large lan-  
 723 guage models. *arXiv preprint arXiv:2505.06304*, 2025.

724

725 Kang Yang, Run Wang, and Lina Wang. Metafinger: Fingerprinting the deep neural networks with  
 726 meta-training. In *Proceedings of the International Joint Conference on Artificial Intelligence*  
 727 (*IJCAI*), 2022.

728

729 Zhiguang Yang and Hanzhou Wu. A fingerprint for large language models. *arXiv preprint*  
 730 *arXiv:2407.01235*, 2024.

731

732 Do-hyeon Yoon, Minsoo Chun, Thomas Allen, Hans Müller, Min Wang, and Rajesh Sharma. In-  
 733 trinsic fingerprint of llms: Continue training is not all you need to steal a model! *arXiv preprint*  
 734 *arXiv:2507.03014*, 2025.

735

736 Boyi Zeng, Lizheng Wang, Yuncong Hu, Yi Xu, Chenghu Zhou, Xinbing Wang, Yu Yu, and Zhouhan  
 737 Lin. Huref: Human-readable fingerprint for large language models. *Advances in Neural Infor-*  
 738 *mation Processing Systems*, 37:126332–126362, 2024.

739

740 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in neural infor-*  
 741 *mation processing systems*, 32, 2019.

742

743 Jie Zhang, Jiayuan Li, Haiqiang Fei, Lun Li, and Hongsong Zhu. Easydetector: Using linear probe  
 744 to detect the provenance of large language models. In *IEEE International Conference on Trust,*  
 745 *Security and Privacy in Computing and Communications*, 2024.

746

747 Jie Zhang, Dongrui Liu, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. REEF:  
 748 Representation encoding fingerprints for large language models. In *The Thirteenth International*  
 749 *Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=SnDmPkOJ0T>.

750

751 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-  
 752 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer  
 753 language models. *arXiv preprint arXiv:2205.01068*, 2022.

754

755 Jingjing Zhao, Qingyue Hu, Gaoyang Liu, Xiaoqiang Ma, Fei Chen, and Mohammad Mehedi Has-  
 756 san. Afa: Adversarial fingerprinting authentication for deep neural networks. *Computer Commu-*  
 757 *nications*, 150:488–497, 2020.

758

759 Xuandong Zhao, Yu-Xiang Wang, and Lei Li. Protecting language generation models via invisible  
 760 watermarking. In *International Conference on Machine Learning*, pp. 42187–42199. PMLR,  
 761 2023.

762

763 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,  
 764 Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. Codegeex: A pre-trained model for  
 765 code generation with multilingual evaluations on humaneval-x. In *KDD*, 2023.

756 **A ETHICS STATEMENT**  
757758 The primary motivation for this work is to protect the intellectual property of large language models  
759 (LLMs), thereby promoting fairness and accountability within the AI community. We acknowl-  
760 edge the potential societal impact of such a technology, particularly the risk of false accusations of  
761 model theft. Therefore, a core design principle of our method is to achieve an exceptionally low  
762 false-positive rate, as empirically demonstrated in our experiments. By establishing a high-fidelity  
763 verification system, we aim to foster a more transparent and trustworthy open-source ecosystem. All  
764 models used in this study are publicly available, and our research does not involve any private or  
765 sensitive user data.766 **B REPRODUCIBILITY STATEMENT**  
767768 To ensure the reproducibility of our results, we have included the complete source code in the sup-  
769 plementary material. We also provided a detailed description of our methodology in [Section 5](#), in-  
770 cluding a step-by-step algorithm. All models used for evaluation are publicly available from sources  
771 such as the Hugging Face Hub, and a comprehensive list mapping our model abbreviations to their  
772 full names is included in [Appendix G](#). A public code repository will be made available upon publi-  
773 cation.  
774775 **C THE USE OF LARGE LANGUAGE MODELS**  
776777 The core research and analysis presented in this manuscript were conducted without the use of Large  
778 Language Models (LLMs). An LLM was utilized exclusively to improve the language and clarity  
779 of the text.  
780781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

810 **D DEFINITIONS AND NOTATIONS**  
 811

812 **D.1 LLM ARCHITECTURE**  
 813

814 **Weights** Here we provide notations for the rest of model A’s weights. In addition to the notations  
 815 in [Section 3](#), we further denote  $s$ ,  $W_{A,\text{lm}}$  as the language model head, and  $W_{A,\text{v}}^{(l)}$ ,  $W_{A,\text{o}}^{(l)}$ ,  $\mathbb{W}_{A,\text{ffn}}^{(l)}$  the  
 816 V,O matrices and FFN at the  $l$ -th layer, and

817 
$$\begin{aligned} 818 \mathbb{W}_{A,\text{others}} &= \mathbb{W}_A / \mathbb{W}_{A,\text{partial}} \\ 819 &= \{W_{A,\text{lm}}\} \bigcup_{l=1}^L \mathbb{W}_{A,\text{ffn}}^{(l)} \bigcup_{l=1}^L \{W_{A,i}^{(l)} \mid i \in \{\text{V, O}\}\} \end{aligned}$$
  
 820

821 as the rest of model A’s weights.

822 **Functions** Here we provide definitions for functions in an LLM.

823 **Definition D.1 (RMSNorm)** *Given the input matrix  $X \in \mathbb{R}^{m \times n}$  of model A, the learnable parameters  $\omega_A^{(l)} = (\omega_{A,1}^{(l)}, \dots, \omega_{A,n}^{(l)}) \in \mathbb{R}^n$  at the  $l$ -th layer, and a constant  $\epsilon_A \in \mathbb{R}$  (this constant is usually shared across all layers of an LLM), RMSNorm at layer  $l$  is a function that*

824 
$$\text{RMSNorm}(X, \omega_A^{(l)}, \epsilon_A) = \text{diag}^{-\frac{1}{2}}(X X^\top \mathbf{1}_m + \epsilon_A \mathbf{1}_m) \cdot X \cdot \text{diag}(\omega_A^{(l)}). \quad (6)$$
  
 825

826 **Definition D.2 (Self-Attention with RoPE)** *Given an input  $X$  to model A, the self-attention mechanism at layer  $l$  is a function  $f_{\text{attn}}$  that*

827 
$$f_{\text{attn}}^{(l)}(X) = \text{AttnScore}(X, W_{A,Q}^{(l)}, W_{A,K}^{(l)}, \theta) \cdot X W_{A,V}^{(l)\top} \cdot W_{A,O}^{(l)\top} \quad (7)$$
  
 828

829 where  $\text{AttnScore}$  is the attention score function with RoPE, i.e.

830 
$$\text{AttnScore}(X, W_{A,Q}^{(l)}, W_{A,K}^{(l)}, \theta) = \text{softmax}\left(\frac{1}{\sqrt{d}} \text{RoPE}(X W_{A,Q}^{(l)\top}, \theta) \cdot \text{RoPE}(X W_{A,K}^{(l)\top}, \theta)^\top\right), \quad (8a)$$
  
 831

832 
$$\text{RoPE}(X W_{A,Q}^{(l)\top}, \theta) = [x_1^\top W_{A,Q}^{(l)\top} R_\theta(0), \dots, x_m^\top W_{A,Q}^{(l)\top} R_\theta(m-1)], \quad (8b)$$
  
 833

834 
$$\text{RoPE}(X W_{A,K}^{(l)\top}, \theta) = [x_1^\top W_{A,K}^{(l)\top} R_\theta(0), \dots, x_m^\top W_{A,K}^{(l)\top} R_\theta(m-1)]. \quad (8c)$$
  
 835

836 Here  $x_i^\top \in \mathbb{R}^n$  is the  $i$ -th row of  $X$ ,  $\theta$  is the frequency base for RoPE and  $R_\theta$  is the rotation  
 837 matrix ([Su et al., 2024](#)).

838 RoPE’s definition in [Definition D.2](#) is slightly different from their implementations, but still an  
 839 equivalent version to [Su et al. \(2024\)](#).

840 **Definition D.3 (Feed-forward Networks)** *Given the input  $X$ , the feed-forward network which  
 841 adopts the SwiGLU MLP is the function  $f_{\text{ffn}}^{(l)}$  defined by*

842 
$$f_{\text{ffn}}^{(l)}(X) = \left( \text{SiLU}(X W_{A,\text{gate}}^{(l)\top}) \odot (X W_{A,\text{up}}^{(l)\top}) \right) W_{A,\text{down}}^{(l)\top},$$
  
 843

844 where  $W_{A,\text{up}}^{(l)}$ ,  $W_{A,\text{gate}}^{(l)}$ ,  $W_{A,\text{down}}^{(l)}$  are bias-free weight matrices,  $\odot$  denotes the Hadamard product,  
 845 and  $\text{SiLU}(z) = z \sigma(z)$  with  $\sigma(z) = \frac{1}{1+e^{-z}}$ . We write  $\mathbb{W}_{A,\text{ffn}}^{(l)} = \{W_{A,\text{up}}^{(l)}, W_{A,\text{gate}}^{(l)}, W_{A,\text{down}}^{(l)}\}$ .  
 846

847 We abuse the notations in [Definition D.1](#) by denoting  $\text{RMSNorm}_{\text{attn}}^{(l)}(X) \triangleq$   
 848  $\text{RMSNorm}(X, \omega_{A,\text{attn}}^{(l)}, \epsilon_A)$  as the attention norm and  $\text{RMSNorm}_{\text{ffn}}^{(l)}(X) \triangleq \text{RMSNorm}(X, \omega_{A,\text{ffn}}^{(l)}, \epsilon_A)$   
 849 the ffn norm at layer  $l$  of model A. Then, the definition of a Transformer block at layer  $l$  of model  
 850 A can be summarized as follows.

851 **Definition D.4 (Transformer Block)** *Combining all the components with residual connections, the  
 852 transformer block at the  $l$ -th layer of model A is a function  $f_{\text{Transformer}}^{(l)}$  that*

853 
$$f_{\text{Transformer}}^{(l)} = \left( I + f_{\text{ffn}}^{(l)} \circ \text{RMSNorm}_{\text{ffn}}^{(l)} \right) \circ \left( I + f_{\text{attn}}^{(l)} \circ \text{RMSNorm}_{\text{attn}}^{(l)} \right). \quad (9)$$
  
 854

864 where  $I$  is the identity mapping,  $f_{\text{attn}}^{(l)}$  and  $f_{\text{ffn}}^{(l)}$  are the attention and FFN at  $l$ -th layer respectively,  
 865  $\text{RMSNorm}_{\text{ffn}}^{(l)}$  and  $\text{RMSNorm}_{\text{attn}}^{(l)}$  are RMSNorm functions for these two components.  
 866

867 **Definition D.5 (LLM Architecture)** Given an word embedding  $X$ , the last hidden state  $Y$  of an  
 868  $L$ -layered LLM  $A$  is  
 869

$$870 \quad Y = f_{\text{Transformer}}^{(L)} \circ f_{\text{Transformer}}^{(L-1)} \circ \cdots \circ f_{\text{Transformer}}^{(1)}(X) \cdot W_{A,lm}^{\top}. \quad (10)$$

## 872 D.2 CENTRAL KERNEL ALIGNMENT (CKA)

874 Given two input matrices  $X_1 \in \mathbb{R}^{m \times n_1}$  and  $X_2 \in \mathbb{R}^{m \times n_2}$ , CKA is a function mapping paired input  
 875 matrices to  $[0, 1]$ , and

$$877 \quad \text{CKA}(X_1, X_2) = \frac{\text{HSIC}(K_{X_1}, K_{X_2})}{\sqrt{\text{HSIC}(K_{X_1}, K_{X_1}) \cdot \text{HSIC}(K_{X_2}, K_{X_2})}} \quad (11)$$

879 where  $\text{HSIC}(K_{X_1}, K_{X_2}) = \frac{1}{(m-1)^2} \text{tr}(K_{X_1} H_m K_{X_2} H_m)$ ,  $(K_{X_1})_{ij} = k((X_1)_i, (X_1)_j)$ ,  $(K_{X_2})_{ij} =$   
 880  $k((X_2)_i, (X_2)_j)$  are kernel matrices with the kernel function  $k(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_+$ , and  $H_m =$   
 881  $I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^{\top}$  is the centering matrix. The linear kernels give  $K_X = X X^{\top}$  for a matrix  $X$ . To  
 882 reduce the finite-sample bias of Eq. (6), we use the unbiased HSIC estimator (Song et al., 2007). Let  
 883  $\tilde{K}_{X_i}$  be  $K_{X_i}$  with its diagonal set to zero, i.e.,  $(\tilde{K}_{X_i})_{ii} = 0$ . With  $\mathbf{1}_m$  the all-ones vector of length  
 884  $m$ , the unbiased HSIC is  
 885

$$886 \quad \text{HSIC}_u(\tilde{K}_{X_1}, \tilde{K}_{X_2}) = \frac{1}{m(m-3)} \left[ \text{tr}(\tilde{K}_{X_1} \tilde{K}_{X_2}) + \frac{(\mathbf{1}_m^{\top} \tilde{K}_{X_1} \mathbf{1}_m)(\mathbf{1}_m^{\top} \tilde{K}_{X_2} \mathbf{1}_m)}{(m-1)(m-2)} - \frac{2}{m-2} \mathbf{1}_m^{\top} \tilde{K}_{X_1} \tilde{K}_{X_2} \mathbf{1}_m \right]. \quad (12)$$

888 We then define

$$890 \quad \text{UCKA}(X_1, X_2) = \frac{\text{HSIC}_u(\tilde{K}_{X_1}, \tilde{K}_{X_2})}{\sqrt{\text{HSIC}_u(\tilde{K}_{X_1}, \tilde{K}_{X_1}) \cdot \text{HSIC}_u(\tilde{K}_{X_2}, \tilde{K}_{X_2})}}, \quad (13)$$

893 which preserves Theorem 3.1 and yields values in  $[-1, 1]$ .

918 **E PROOFS AND DISCUSSIONS**  
 919

920 **E.1 PROOF FOR THEOREM 3.1**  
 921

922 We give the proof based on the definition of CKA with linear kernels in [Appendix D.2](#). Given input  
 923 matrices  $X_1, X_2$ , orthogonal matrices  $U_1, U_2$  and constants  $c_1, c_2$ , the kernel functions yield  
 924

$$925 \quad K_{c_1 X_1 U_1} = c_1 X_1 U_1 U_1^\top X_1^\top c_1 = c_1^2 X_1 X_1^\top = c_1^2 K_{X_1}, \\ 926 \quad K_{c_2 X_2 U_2} = c_2 X_2 U_2 U_2^\top X_2^\top c_2 = c_2^2 X_2 X_2^\top = c_2^2 K_{X_2}.$$

928 Therefore, corresponding HSIC results are  
 929

$$930 \quad \text{HSIC}(K_{c_1 X_1 U_1}, K_{c_2 X_2 U_2}) = \frac{c_1^2 c_2^2}{(m-1)^2} \text{tr}(K_{X_1} H_m K_{X_2} H_m) = c_1^2 c_2^2 \text{HSIC}(K_{X_1}, K_{X_2}), \\ 931 \\ 932 \quad \text{HSIC}(K_{c_1 X_1 U_1}, K_{c_1 X_1 U_1}) = \frac{c_1^4}{(m-1)^2} \text{tr}(K_{X_1} H_m K_{X_1} H_m) = c_1^4 \text{HSIC}(K_{X_1}, K_{X_1}), \\ 933 \\ 934 \quad \text{HSIC}(K_{c_2 X_2 U_2}, K_{c_2 X_2 U_2}) = \frac{c_2^4}{(m-1)^2} \text{tr}(K_{X_2} H_m K_{X_2} H_m) = c_2^4 \text{HSIC}(K_{X_2}, K_{X_2}).$$

937 Hence, it follows that  
 938

$$939 \quad \text{CKA}(c_1 X_1 U_1, c_2 X_2 U_2) = \frac{\text{HSIC}(K_{c_1 X_1 U_1}, K_{c_2 X_2 U_2})}{\sqrt{\text{HSIC}(K_{c_1 X_1 U_1}, K_{c_1 X_1 U_1}) \cdot \text{HSIC}(K_{c_2 X_2 U_2}, K_{c_2 X_2 U_2})}} \\ 940 \\ 941 \quad = \frac{c_1^2 c_2^2 \text{HSIC}(K_{X_1}, K_{X_2})}{c_1^2 c_2^2 \sqrt{\text{HSIC}(K_{X_1}, K_{X_1}) \cdot \text{HSIC}(K_{X_2}, K_{X_2})}} \\ 942 \\ 943 \quad = \frac{\text{HSIC}(K_{X_1}, K_{X_2})}{\sqrt{\text{HSIC}(K_{X_1}, K_{X_1}) \cdot \text{HSIC}(K_{X_2}, K_{X_2})}} \\ 944 \\ 945 \quad = \text{CKA}(X_1, X_2).$$

948 **E.2 PROOF FOR PROPOSITION 4.2**  
 949

950 The goal is to show the manipulation  $T$  commutes with  $f$ , i.e.  $f \circ T = T \circ f$ . If  $f = f_n \circ \dots \circ f_1$ ,  
 951 then it suffices to show  
 952

$$953 \quad f_n \circ \dots \circ f_1 \circ T = T \circ f_n \circ \dots \circ f_1.$$

954 Since  
 955

$$956 \quad f_1 \circ T = T \circ f_1, f_2 \circ T = T \circ f_2, \dots, f_n \circ T = T \circ f_n,$$

957 It is direct that  
 958

$$959 \quad f \circ T = T \circ f.$$

960 **E.3 PROOF FOR THEOREM 4.3 AND DISCUSSIONS ON OTHER INPUT MANIPULATIONS**  
 961

962 **We first show a proof for the claims in Theorem 4.3.** The core lies in preserving the diagonal  
 963 structure of RMSNorm parameters in [Definition D.1](#). Given the input  $X$  to layer  $l$ , the manipulations  
 964 change it to  $cXPD$  and fed the manipulated input into the attention norm of model B. In order to  
 965 recover the manipulations on inputs after RMSNorm, i.e. for any manipulation  $T$  in the input  $X$  (an  
 966 example is  $T(X) = cXPD$ ),  
 967

$$968 \quad \text{RMSNorm}_{\text{attn}}^{(l)} \circ T = T \circ \text{RMSNorm}_{\text{attn}}^{(l)}, \quad (14)$$

969 it suffices to show  
 970

$$971 \quad \exists \omega_{B, \text{attn}}^{(l)} \text{ and } \epsilon_B, \text{ s.t. } \text{RMSNorm}(X, \omega_{A, \text{attn}}^{(l)}, \epsilon_A) \cdot cPD = \text{RMSNorm}(cXPD, \omega_{B, \text{attn}}^{(l)}, \epsilon_B)$$

since the learnable parameters and norm epsilon can be modified to satisfy [Equation 14](#) and [Proposition 4.2](#). By setting  $\epsilon_B = c^2\epsilon_A$  and  $\text{diag}(\omega_{B,\text{attn}}^{(l)}) = cD^\top P^\top \text{diag}(\omega_{A,\text{attn}}^{(l)})PD$ , one has

$$\begin{aligned} & \text{RMSNorm}(cXPD, \omega_{B,\text{attn}}^{(l)}, \epsilon_B) \\ &= \text{diag}^{-\frac{1}{2}}((cXPD)(cXPD)^\top \mathbf{1}_m + \epsilon_B \mathbf{1}_m) \cdot cXPD \cdot \text{diag}(\omega_{B,\text{attn}}^{(l)}) \\ &= \text{diag}^{-\frac{1}{2}}(c^2 XX^\top \mathbf{1}_m + c^2 \epsilon_A \mathbf{1}_m) \cdot cXPD \cdot cD^\top P^\top \text{diag}(\omega_{A,\text{attn}}^{(l)})PD \\ &= c \cdot \text{diag}^{-\frac{1}{2}}(XX^\top \mathbf{1}_m + \epsilon_A \mathbf{1}_m) \cdot X \cdot \text{diag}(\omega_{A,\text{attn}}^{(l)}) \cdot PD \\ &= \text{RMSNorm}(X, \omega_{A,\text{attn}}^{(l)}) \cdot cPD. \end{aligned}$$

The proof still holds for partial permutations  $P$ , i.e.  $P$  is rectangular with number of rows more than number of columns, since  $cD^\top P^\top \text{diag}(\omega_{A,\text{attn}}^{(l)})PD$  is still a diagonal matrix. Therefore,  $R_{B,\text{emb}} = cPD$  under [Proposition 4.2](#). **However, for more general manipulations, Proposition 4.2 does not hold for RMSNorm.** We give two examples for illustration.

First, **orthogonal transformations on the input generally fail to pass through RMSNorm.** Given any orthogonal matrices  $U$  as the manipulation on  $X$ , i.e.,  $T(X) = XU$ , if one attempts to keep [Equation 14](#), then it requires

$$\exists \omega^{(l)} \text{ and } \epsilon, \text{ s.t. } \text{RMSNorm}(XU, \omega^{(l)}, \epsilon) = \text{RMSNorm}(X, \omega_{A,\text{attn}}^{(l)}, \epsilon_A)U.$$

However,

$$\text{RMSNorm}(XU, \omega, \epsilon) = \text{diag}^{-\frac{1}{2}}(XX^\top \mathbf{1}_m + \epsilon \mathbf{1}_m) \cdot XU \cdot \text{diag}(\omega^{(l)}).$$

Hence, it is required that

$$U \cdot \text{diag}(\omega^{(l)}) = \text{diag}(\omega_{A,\text{attn}}^{(l)})U$$

which suggests  $U^\top \text{diag}(\omega_{A,\text{attn}}^{(l)})U$  is a diagonal matrix. Nevertheless, this property generally does not hold if  $U$  is not a (partial) permutation matrix  $P$  or a signature matrix  $D$ .

**Second, non-orthogonal transformations on inputs generally fail to pass through RMSNorm.** Given an arbitrary (invertible) transformation  $M$ , [Proposition 4.2](#) and [Equation 14](#) require that

$$\exists \omega^{(l)} \text{ and } \epsilon, \text{ s.t. } \text{RMSNorm}(XM, \omega^{(l)}, \epsilon) = \text{RMSNorm}(X, \omega_{A,\text{attn}}^{(l)}, \epsilon_A)M.$$

However, since

$$\text{RMSNorm}(XM, \omega^{(l)}, \epsilon) = \text{diag}^{-\frac{1}{2}}(XMM^\top X^\top \mathbf{1}_m + \epsilon \mathbf{1}_m) \cdot XM \cdot \text{diag}(\omega^{(l)}),$$

any  $M$  that do not satisfy  $MM^\top = c'I$  where  $c'$  is a constant and  $I$  is the identity matrix can hardly recover the manipulations due to the nonlinearity in norm functions. Moreover, similar to the case with orthogonal transformations,  $M \cdot \text{diag}(\omega^{(l)}) = \text{diag}(\omega_{A,\text{attn}}^{(l)})M$  also requires  $M$  to be (partial) permutation or signature matrices (or a combination of the both).

**Third, we additionally clarify that although the combination of (partial) permutation and signature matrices can be in any order, it is reasonable to fix it as  $R_{B,\text{emb}} = cPD$ .** To be specific, we clarify that for any (partial) permutation matrix  $P'$  and signature matrix  $D'$ , there exists a (partial) permutation matrix  $P$  and a signature  $D$  such that

$$D'P' = PD$$

The proof is straightforward. Define  $P$  entry-wise by

$$P_{ij} := |(D'P')_{ij}| \in \{0, 1\}.$$

Since left-multiplication by  $D'$  only flips signs, each column of  $D'P'$  has at most one nonzero entry. Hence,  $P$  is a partial permutation matrix. For each column  $j$ , if that column of  $D'P'$  has its unique nonzero at row  $i$ , set

$$D_{jj} := \text{sgn}((D'P')_{ij}) \in \{\pm 1\}.$$

If the column is entirely zero, choose  $D_{jj} \in \{\pm 1\}$  arbitrarily. Then for all  $i, j$ ,

$$(PD)_{ij} = P_{ij}D_{jj} = |(D'P')_{ij}| \cdot \text{sgn}((D'P')_{ij}) = (D'P')_{ij},$$

1026 where we interpret  $\text{sgn}(0) = 1$ . Therefore,  $PD = D'P'$ .  
 1027

1028 **Last, we clarify that the conclusions in Theorem 4.3 can generalize to the case where model A**  
 1029 **and model B may not share the same architecture.** It is a direct result from the fact that the two  
 1030 LLMs have shared tokens and that the pruning over dimensions of a model’s word embeddings can  
 1031 be viewed as a partial permutation matrix.

1032 **E.4 PROOF FOR THEOREM 4.4 AND RELATED DISCUSSIONS**  
 1033

1034 **We first provide the proof for Theorem 4.4 based on Definition D.2.**  
 1035

1036 We denote by  $X$  the original input to self-attention, and  $X'$  the manipulated one. Then, Theorem 4.3  
 1037 suggests that

$$1038 \quad X' = cXPD$$

1039 Therefore, the attention score of model B becomes

$$1040 \quad \text{AttnScore}(cXPD, W_{B,Q}^{(l)}, W_{B,K}^{(l)}, \theta) = \text{softmax}\left(\frac{1}{\sqrt{d}}\text{RoPE}(cXPD W_{B,Q}^{(l)\top}, \theta) \cdot \text{RoPE}(cXPD W_{B,K}^{(l)\top}, \theta)^\top\right).$$

1042 To keep attention scores of model B identical to that of model A, it is required that  
 1043

$$1044 \quad \text{RoPE}(cXPD W_{B,Q}^{(l)\top}, \theta) \cdot \text{RoPE}(cXPD W_{B,K}^{(l)\top}, \theta)^\top = \text{RoPE}(X W_{A,Q}^{(l)\top}, \theta) \cdot \text{RoPE}(X W_{A,K}^{(l)\top})^\top.$$

1046 By Equation 8b and Equation 8c, this translates into

$$1047 \quad c^2 x_i^\top P D W_{B,Q}^{(l)\top} R_\theta(i) R_\theta^\top(j) W_{B,K}^{(l)} D^\top P^\top x_j = x_i^\top W_{A,Q}^{(l)\top} R_\theta(i) R_\theta^\top(j) W_{A,K}^{(l)} x_j$$

1049 Therefore, the recovery of model A’s attention score requires column-wise transformations to eliminate  
 1050 the **input-related manipulations passed by RMSNorm**, i.e.

$$1051 \quad W_{B,Q}^{(l)} = c^{-1} W_{A,Q}^{(l)} P D, \quad W_{B,K}^{(l)} = c^{-1} W_{A,K}^{(l)} P D \quad (15)$$

1053 which suggests

$$1054 \quad R_{B,Q}^{(l)} = R_{B,K}^{(l)} = c^{-1} P D.$$

1056 Although Equation 15 recovers attention scores, the Q,K matrices can still be modified without  
 1057 changing attention scores. An example is, given a rotation matrix  $R = R_\theta(k)$  with the same fre-  
 1058 quency base as RoPE, multiplying Q,K matrices with it does not change the attention scores, i.e.

$$1059 \quad x_i^\top W_{B,Q}^{(l)\top} R_\theta(k) R_\theta(i) R_\theta^\top(j) R_\theta^\top(k) W_{B,K}^{(l)} x_j = x_i^\top W_{B,Q}^{(l)\top} R_\theta(i+k) R_\theta^\top(j+k) W_{B,K}^{(l)\top} x_j$$

$$1060 \quad = x_i^\top W_{B,Q}^{(l)\top} R_\theta(i-j) W_{B,K}^{(l)} x_j$$

$$1061 \quad = x_i^\top W_{B,Q}^{(l)\top} R_\theta(i) R_\theta^\top(j) W_{B,K}^{(l)\top} x_j.$$

1064 Furthermore, this conclusion can be generalize to any rotation matrix with structures similar to  
 1065 RoPE rotation matrices. Given  $R = \text{diag}(R(\psi_0), R(\psi_1), \dots)$  with  $R(\psi_k) = \begin{bmatrix} \cos \psi_k & -\sin \psi_k \\ \sin \psi_k & \cos \psi_k \end{bmatrix}$ ,  
 1066 multiplying Q,K matrices by  $R$  keeps the attention scores since  
 1067

$$1068 \quad R R_\theta(i) R_\theta^\top(j) R^\top = R_\theta(i) R_\theta^\top(j). \quad (16)$$

1070 Hence, there are various rotation matrices to significantly change Q,K matrices but preserve attention  
 1071 scores. Since these rotation matrices are naturally orthogonal, we reformulate the row-wise  
 1072 transformations on Q,K matrices  $L_{B,Q}, L_{B,K}$  as  $U_{B,Q}^{(l)}$  and  $U_{B,K}^{(l)}$ .  
 1073

1074 Therefore, we denote **RoPE-related manipulations as**

$$1075 \quad W_{B,Q}^{(l)} = U_{A,Q}^{(l)} W_{B,Q}^{(l)}, \quad W_{B,K}^{(l)} = U_{A,K}^{(l)} W_{B,K}^{(l)}.$$

1077 On the other hand, semi-orthogonal  $U_{B,Q}^{(l)}, U_{B,K}^{(l)}$  are also possible. This is because of the common  
 1078 practice in pruning may select different rows of Q,K matrices, which suggests multiplying Q,K  
 1079 matrices with row-wise partial permutations.

1080  
1081 **Next, we show how Proposition 4.2 is satisfied at the self-attention mechanism based on Theorem 4.4.** Let  $W_{B,V}^{(l)} = W_{A,V}^{(l)}PD$  and  $W_{B,O}^{(l)} = D^\top P^\top W_{A,O}^{(l)}$ . Then,

$$1083 \quad X'W_{B,V}^{(l)\top}W_{B,O}^{(l)\top} = cXPD \cdot D^\top P^\top W_{A,V}^{(l)\top}W_{A,O}^{(l)\top}PD \\ 1084 \\ 1085 \quad = XW_{A,V}^{(l)\top}W_{A,O}^{(l)\top} \cdot cPD$$

1086 recovers the manipulation over inputs in [Equation 7](#).

1088 **E.5 HOW POTENTIAL ATTACKS IN [SECTION 4.5](#) RECOVERS THE OUTPUT OF**  
1089 **TRANSFORMERS**

1091 We provide an illustration based on [Definition D.3](#), [Definition D.4](#) and [Definition D.5](#). Previous parts  
1092 have demonstrated how the manipulations on inputs pass through self-attention. Hence, it suffices  
1093 to show how the manipulations pass through FFN and yield an output identical to model A after the  
1094 language model head.

1095 For layer  $l$ , let  $W_{B,\text{gate}}^{(l)} = c^{-1}W_{A,\text{gate}}^{(l)}PD$ ,  $W_{B,\text{up}}^{(l)} = c^{-1}W_{A,\text{up}}^{(l)}PD$ ,  $W_{B,\text{down}}^{(l)} = cD^\top P^\top W_{A,\text{down}}^{(l)}$ .  
1096 We denote  $X' = cXPD$  as the manipulated input. Then for model B,

$$1097 \quad f_{\text{ffn}}^{(l)}(X') \\ 1098 \quad = \left( \text{SiLU}(X'W_{B,\text{gate}}^{(l)\top}) \odot (X'W_{B,\text{up}}^{(l)\top}) \right) W_{B,\text{down}}^{(l)\top} \\ 1099 \\ 1100 \quad = \left( \text{SiLU}(cXPD \cdot D^\top P^\top W_{A,\text{gate}}^{(l)\top} \cdot c^{-1}) \odot (cXPD \cdot D^\top P^\top W_{A,\text{up}}^{(l)\top} \cdot c^{-1}) \right) W_{B,\text{down}}^{(l)\top} \cdot cPD \\ 1101 \\ 1102 \quad = \left( \text{SiLU}(XW_{A,\text{gate}}^{(l)\top}) \odot (XW_{A,\text{up}}^{(l)\top}) \right) W_{A,\text{down}}^{(l)\top} \cdot cPD.$$

1103 By [Definition D.4](#) and [Definition D.5](#), the manipulation is propagated through transformer layers.  
1104 At the language model head, we denote  $W_{B,\text{lm}} = c^{-1}W_{A,\text{lm}}PD$  and abuse  $X' = cXPD$  as the  
1105 manipulated input. Then,

$$1106 \quad X'W_{B,\text{lm}}^{(l)\top} = cXPD \cdot D^\top P^\top W_{A,\text{lm}}^{(l)\top} \cdot c^{-1} \\ 1107 \\ 1108 \quad = XW_{A,\text{lm}}^{(l)\top}$$

1109 recovers the output of model A.

## F ABLATION STUDIES

## F.1 NUMBER OF OVERLAPPING VOCABULARY TOKENS

Although Algorithm 1 uses overlapping tokens to recover the signature matrices and permutations, the detection performance of AWM does not heavily depend on the amount of vocabulary overlap. In fact, AWM remains effective even when only a small number of tokens (100 tokens) are shared between the two vocabularies. To quantify this, we conduct an ablation study on the number of overlapping tokens used in Algorithm 1. Specifically, for each scenario in Table 2 (SFT, Continual Pretraining, Upcycling, Multi-Modal, RL, and Pruning), we compute the average absolute Z-score under different numbers of overlapping vocabulary tokens and report the results in Figure 3.

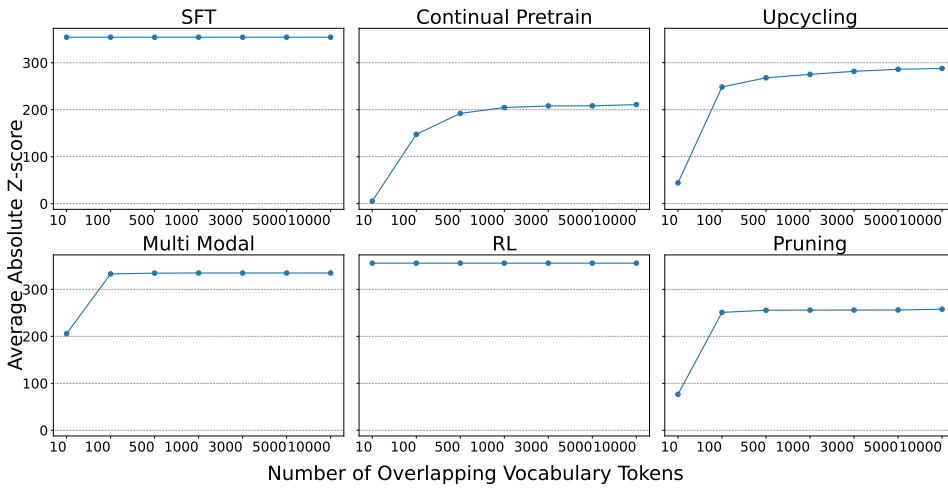


Figure 3: Ablation studies on the number of overlapping vocabulary tokens vs AWM’s average absolute Z-score in Table 2. AWM remain effective even when there are only 100 overlapping tokens used in Algorithm 1.

## F.2 CKA ABLATIONS

We further conduct an ablation study on the design of CKA in Algorithm 1. In AWM, we use the unbiased version with linear kernel for computation efficiency and accuracy. Here we investigate into more variants of CKA, and summarize the results in Table 4. It can be seen that the choice of unbiasedness is crucial to the robustness of our method. Meanwhile, although RBF kernels show stronger performance in Table 4, the choice of linear kernel has already yield strong performance, which we choose for better computational efficiency.

| CKA Kernel         | Linear         |          | RBF      |          |        |
|--------------------|----------------|----------|----------|----------|--------|
|                    | CKA Biasedness | Unbiased | Biased   | Unbiased | Biased |
| SFT                | 356.0223       | 18.3099  | 527.3949 | 11.5843  |        |
| Continual Pretrain | 217.5003       | 11.2954  | 320.9173 | 7.0701   |        |
| Upcycling          | 291.6191       | 15.1994  | 432.1476 | 9.6776   |        |
| Multi Modal        | 336.6757       | 17.3766  | 498.8918 | 11.0056  |        |
| RL                 | 357.5001       | 18.3850  | 529.6448 | 11.6330  |        |
| Pruning            | 268.9175       | 13.8985  | 394.8391 | 8.7331   |        |

Table 4: AWM’s average absolute Z-scores under CKA with different kernels and biasedness.

1188 **G MODEL DETAILS**

1190 **Table 5** details the specifications of the offspring models analyzed in our study, mapping abbreviations to their base models and training datasets. To facilitate reproducibility, each entry in the "Full  
1191 Model Name" column serves as a direct link to the official model checkpoint hosted on the Hugging  
1192 Face Hub.  
1193

1194 **Table 5: Mapping of model abbreviations to their full model names, corresponding Hugging Face  
1195 checkpoints url, base models, and relevant training corpus information.**  
1196

| Abbreviation | Base Model   | Full Model Name                                       | Train Corpus                               |
|--------------|--------------|-------------------------------------------------------|--------------------------------------------|
| Vicuna       | Llama2-7B    | <a href="#">vicuna-7b-v1.5</a>                        | ShareGPT                                   |
| Selfrag      | Llama2-7B    | <a href="#">selfrag_llama2_7b</a>                     | Self-RAG Data                              |
| 32K          | Llama2-7B    | <a href="#">LLaMA-2-7B-32K</a>                        | Book, ArXiv, etc.                          |
| Wizard       | Llama2-7B    | <a href="#">WizardMath-7B-V1.0</a>                    | GSM8k                                      |
| Guanaco      | Llama2-13B   | <a href="#">llama-2-7b-guanaco</a>                    | OpenAssistant                              |
| Vicuna       | Llama2-13B   | <a href="#">vicuna-13b-v1.5</a>                       | ShareGPT                                   |
| Hermes       | Llama2-13B   | <a href="#">Nous-Hermes-Llama2-13b</a>                | GPTeacher, WizardLM, etc.                  |
| Estopia      | Llama2-13B   | <a href="#">LLaMA2-13B-Estopia</a>                    | EstopiaV9/V13, Tiefighter, etc.            |
| Finance      | Llama2-7B    | <a href="#">llama-2-7b-finance</a>                    | Financial Dataset                          |
| Firefly      | Llama2-13B   | <a href="#">firefly-llama2-13b</a>                    | CLUE, ThucNews, etc.                       |
| Llemma       | Llama2-7B    | <a href="#">llemma_7b</a>                             | ArXiv, OpenWebMath, etc.                   |
| Code         | Llama2-7B    | <a href="#">CodeLlama-7b-hf</a>                       | Deduped code, Natural language             |
| Python       | Llama2-7B    | <a href="#">CodeLlama-7b-Python-hf</a>                | Python code                                |
| Code         | Gemma-2B     | <a href="#">codegemma-2b</a>                          | Math, Synthetic code, etc.                 |
| Code         | Gemma-7B     | <a href="#">codegemma-7b</a>                          | Code, Natural language                     |
| Math         | Qwen2.5-7B   | <a href="#">Qwen2.5-Math-7B</a>                       | Web, Books, etc.                           |
| Coder        | Qwen2.5-7B   | <a href="#">Qwen2.5-Coder-7B</a>                      | Source Code, Synthetic data, etc.          |
| Math         | Qwen2-7B     | <a href="#">Qwen2-Math-7B</a>                         | Math Data                                  |
| Code         | Llama2-70B   | <a href="#">CodeLlama-70b-hf</a>                      | Deduped code, Natural language             |
| Python       | Llama2-70B   | <a href="#">CodeLlama-70b-Python-hf</a>               | Python code                                |
| Mixtral      | Mistral-7B   | <a href="#">Nous-Hermes-2-Mixtral-8x7B-DPO</a>        | GPT-4 Data, Open datasets                  |
| MoE v2       | Llama3-8B    | <a href="#">LLaMA-MoE-v2-3_8B-2_8-sft</a>             | SFT Data                                   |
| MoE4         | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_5B-4_16</a>                | SlimPajama                                 |
| MoE 3B       | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_0B-2_16</a>                | SlimPajama                                 |
| MoE2         | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_5B-2_8</a>                 | SlimPajama                                 |
| MoE3B-SFT    | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_0B-2_16-sft</a>            | SlimPajama, SFT Data                       |
| MoE2-SFT     | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_5B-2_8-sft</a>             | SlimPajama, SFT Data                       |
| MoE4-SFT     | Llama2-7B    | <a href="#">LLaMA-MoE-v1-3_5B-4_16-sft</a>            | SlimPajama, SFT Data                       |
| Qwen1.5 MoE  | Qwen-1.8B    | <a href="#">Qwen1.5-MoE-A2.7B</a>                     | Qwen Base Corpus                           |
| Minicpm MoE  | Minicpm-2B   | <a href="#">MiniCPM-MoE-8x2B</a>                      | MiniCPM Data                               |
| LLaVA        | Llama2-7B    | <a href="#">llava-v1.5-7b</a>                         | LAION, GPT instructions, etc.              |
| Video        | Llama2-7B    | <a href="#">Video-LLaVA-7B-hf</a>                     | Caption, QA                                |
| VL           | Qwen2-7B     | <a href="#">Qwen2-VL-7B-Instruct</a>                  | Image-text, OCR, etc.                      |
| Audio        | Qwen-7B      | <a href="#">Qwen-Audio</a>                            | Speech, Sound, etc.                        |
| Audio2       | Qwen-7B      | <a href="#">Qwen2-Audio-7B</a>                        | Audio-text, Voice Chat                     |
| VL           | Qwen-7B      | <a href="#">Qwen-VL</a>                               | Image-text, OCR, etc.                      |
| VL           | Qwen2.5-7B   | <a href="#">Qwen2.5-VL-7B-Instruct</a>                | Visual recognition, Document parsing, etc. |
| VL           | Qwen2.5-3B   | <a href="#">Qwen2.5-VL-3B-Instruct</a>                | Visual recognition, Document parsing, etc. |
| Next         | Llama3-8B    | <a href="#">llama3-llava-next-8b-hf</a>               | LLaVA-NeXT Data                            |
| LLaVA        | Llama2-13B   | <a href="#">llava-v1.5-13b</a>                        | LAION, GPT instructions, etc.              |
| RLHF         | Open-llama3B | <a href="#">hh_rlfh_rm_open_llama_3b</a>              | Anthropic HH-RLHF                          |
| Reason       | Qwen2.5-7B   | <a href="#">Nemotron-Research-Reasoning-Qwen-1.5B</a> | Math, Code, etc.                           |
| Zero         | Qwen2.5-1.5B | <a href="#">Open-Reasoner-Zero-1.5B</a>               | AIME 2024, MATH500, etc.                   |
| DPO          | Mixtral      | <a href="#">Nous-Hermes-2-Mixtral-8x7B-DPO</a>        | GPT-4 Data, Preference pairs               |
| DPO          | Mistral-7B   | <a href="#">Nous-Hermes-2-Mistral-7B-DPO</a>          | GPT-4 Data, Preference pairs               |

1242 **Table 5 – continued from previous page**  
1243

| 1244 <b>Abbreviation</b> | 1245 <b>Base Model</b> | 1246 <b>Full Model Name</b>      | 1247 <b>Train Corpus</b>       |
|--------------------------|------------------------|----------------------------------|--------------------------------|
| Dolphin                  | Mistral-7B             | dolphin-2.6-mistral-7b-dpo       | UltraFeedback, Magicoder, etc. |
| DPO                      | Minicpm-2B             | MiniCPM-2B-dpo-bf16              | ShareGPT, UltraChat, etc.      |
| GRPO                     | Qwen3-4B               | Qwen3_Medical_GRPO               | Medical dataset                |
| RLHF                     | Chatglm-6B             | chatglm-fitness-RLHF             | SFT, Reward Model, etc.        |
| DPO                      | Llama3-8B              | LLAMA3-iterative-DPO-final       | UltraFeedback, Preference sets |
| Minitron-Depth           | Llama-3-8B             | Llama-3.1-Minitron-4B-Depth-Base | Nemotron-4 15B corpus          |
| Minitron-Width           | Llama-3-8B             | Llama-3.1-Minitron-4B-Width-Base | Nemotron-4 15B corpus          |
| Sheared 2.7B-P           | Llama2-7B              | Sheared-LLaMA-2.7B-Pruned        | RedPajama                      |
| Sheared 2.7B-S           | Llama2-7B              | Sheared-LLaMA-2.7B-ShareGPT      | ShareGPT                       |
| Sheared 2.7B             | Llama2-7B              | Sheared-LLaMA-2.7B               | RedPajama                      |
| Sheared 1.3B-P           | Llama2-7B              | Sheared-LLaMA-1.3B-Pruned        | RedPajama                      |
| Sheared 1.3B             | Llama2-7B              | Sheared-LLaMA-1.3B               | RedPajama                      |
| Sheared 1.3B-S           | Llama2-7B              | Sheared-LLaMA-1.3B-ShareGPT      | ShareGPT                       |
| Llama3-1B                | Llama3-8B              | Llama-3.2-1B                     | Llama3-8B logits, Safety data  |
| Llama3-3B                | Llama3-8B              | Llama-3.2-3B                     | Llama3-8B logits, Safety data  |

1260  
1261 

## H IMPLEMENTATION DETAILS

  
1262

1263 We employ the Linear Kernel ( $k(X, Y) = XY^\top$ ) for Centered Kernel Alignment (CKA) due to its  
 1264 computational efficiency. To mitigate the finite-sample bias inherent in standard HSIC estimations,  
 1265 we utilize the Unbiased CKA (UCKA) estimator. As for module selection, our method operates  
 1266 on two specific sets of weights: first, we utilize the intersection of the word embeddings ( $W_{emb}$ )  
 1267 to solve the Linear Assignment Problem (LAP), allowing us to accurately recover the permutation  
 1268 ( $P$ ) and signature ( $D$ ) matrices; second, we compute the final fingerprinting scores using the Query  
 1269 ( $W_Q$ ) and Key ( $W_K$ ) weights, as their transformations are strictly constrained.

1270 To address structural discrepancies such as differing layer counts, we identify the optimal layer  
 1271 correspondence by maximizing the total similarity; specifically, we solve the assignment problem  
 1272 on a cost matrix constructed from the pairwise UCKA scores of  $W_Q$  and  $W_K$  between all source  
 1273 and target layers.

1274  
1275 

## I EMPIRICAL VALIDATION OF ROBUSTNESS AGAINST WEIGHT 1276 MANIPULATIONS

  
1277

1278 In Section 4, we theoretically analyze potential weight manipulations, including constant scaling,  
 1279 signature matrix multiplication, permutations, and orthogonal transformations. We now provide  
 1280 empirical evidence to support the analysis. Specifically, we apply these manipulations to five repre-  
 1281 sentative models, Llama-2-7B, Qwen2-7B, Mistral-7B, Llama-3-8B, and Gemma-7B, and evaluate  
 1282 AWM’s robustness under each setting.

1283 As shown in Table 6, AWM achieves a 100% detection rate across all tested model–manipulation  
 1284 pairs. These results not only align with our theoretical derivations, but also validate the design of  
 1285 AWM, including LAP and UCKA.

1286  
1287 

Table 6: AWM-detected similarity scores under the weight manipulations in Section 4.

  
1288

| 1289 <b>Manipulation / Model</b> | 1290 <b>Llama-2-7B</b> | 1291 <b>Qwen2-7B</b> | 1292 <b>Mistral-7B</b> | 1293 <b>Llama-3-8B</b> | 1294 <b>Gemma-7B</b> |
|----------------------------------|------------------------|----------------------|------------------------|------------------------|----------------------|
| Permutation                      | 100%                   | 100%                 | 100%                   | 100%                   | 100%                 |
| Signature                        | 100%                   | 100%                 | 100%                   | 100%                   | 100%                 |
| Constant Scaling                 | 100%                   | 100%                 | 100%                   | 100%                   | 100%                 |
| Orthogonal Trans.                | 100%                   | 100%                 | 100%                   | 100%                   | 100%                 |