

# Robust Semantic Watermarking for Mixture-of-Experts Large Language Models: A Novel Architecture-Aware Framework

Your N. Here  
Your Institution

Second Name  
Second Institution

## Abstract

Mixture-of-Experts (MoE) models represent a paradigm shift toward sparse computation, yet existing watermarking techniques fail to leverage their unique routing patterns. We introduce MoE-native watermarking that exploits discrete expert selection rather than continuous routing weights. Our framework presents three methods: (1) *Combinatorial Expert Signatures (CES)* with error-correcting codes for deterministic robustness, (2) *Trajectory Graph Hashing (TGH)* for hierarchical semantic encoding, and (3) *Keyed Learnable Quantizer (KLQ)* using contrastive learning for semantic invariance. Experiments demonstrate superior robustness against paraphrase attacks while maintaining competitive efficiency.

## 1 Introduction

Large Language Models (LLMs) based on Mixture-of-Experts (MoE) architectures are becoming prevalent [?], yet existing watermarking methods treat them as dense networks, missing opportunities for robust semantic watermarking. Current approaches like Kirchenbauer et al. [?] rely on token-based hashing, vulnerable to paraphrase attacks.

**Key Insight:** MoE models activate only top- $k$  experts per input, creating discrete routing patterns that are inherently more robust to semantic perturbations than continuous embeddings. We propose *MoE-native watermarking* that exploits this discrete, combinatorial structure.

**Contributions:** We introduce three complementary methods that leverage different aspects of MoE computation, providing deterministic robustness guarantees and superior performance against paraphrase attacks.

## 2 Core Methods

### 2.1 Combinatorial Expert Signatures (CES)

CES encodes the discrete expert selection at each MoE layer using error-correcting codes for algebraic robustness guarantees.

**Core Mechanism:** For input context  $x$ , we extract the top- $k$  expert set:

$$\text{TopK}(x) = \{i_1, i_2, \dots, i_k\} \subset \{1, \dots, E\} \quad (1)$$

We encode this set as a preliminary signature  $s' \in \{0, 1\}^{L_{\text{msg}}}$ :

$$s' = \text{Enc}_{\text{comb}}(\text{TopK}(x)) \quad (2)$$

Then apply error-correcting code with generator matrix  $G$ :

$$s = s' \cdot G \in \{0, 1\}^{L_{\text{code}}} \quad (3)$$

**Key Innovation:** We model paraphrase attacks as communication channel noise. If an attack changes at most  $t$  experts, the resulting signature can be corrected using ECC decoding, providing deterministic robustness guarantees.

### 2.2 Trajectory Graph Hashing (TGH)

TGH captures hierarchical semantic processing by encoding expert activation sequences across multiple MoE layers as graph structures.

**Expert Trajectory:** For input  $x$ , we extract the expert selection sequence:

$$T(x) = (\text{TopK}_1(x), \text{TopK}_2(x), \dots, \text{TopK}_{L_{\text{moe}}}(x)) \quad (4)$$

**Graph-Based Encoding:** We represent the trajectory as a graph and extract structural features:

$$v = \Phi(T(x)) \in \mathbb{R}^D \quad (5)$$

The feature vector is then hashed to produce the final signature:

$$s = H(v) \in \{0, 1\}^L \quad (6)$$

### 2.3 Keyed Learnable Quantizer (KLQ)

KLQ employs contrastive learning to automatically discover semantic-invariant mappings from routing weights to discrete signatures.

**Contrastive Learning:** We train a small auxiliary network  $Q_k$  (parameterized by secret key  $k$ ) using paraphrase pairs as positive examples. The contrastive loss encourages similar outputs for paraphrases:

$$\begin{aligned}\mathcal{L} &= -\log \frac{\exp(s_+/\tau)}{\exp(s_+/\tau) + \sum_j \exp(s_{-j}/\tau)} \\ &= -\log \frac{\exp(s_+/\tau)}{Z(x)}\end{aligned}\quad (7)$$

where  $s_+ = \text{sim}(Q_k(R(x)), Q_k(R(x')))$  and  $s_{-j} = \text{sim}(Q_k(R(x)), Q_k(R(x'_j)))$ .

**Zero-Cost Training:** Only the auxiliary quantizer is trained; the main LLM remains frozen.

### 3 Theoretical Analysis

#### 3.1 Robustness Guarantees

**CES:** For an ECC with minimum distance  $d_{\min}$ , the system can correct up to:

$$t = \lfloor (d_{\min} - 1)/2 \rfloor \quad (8)$$

expert substitutions, providing deterministic robustness guarantees.

**TGH:** The hierarchical structure provides natural robustness through layer-wise semantic abstraction.

**KLQ:** Statistical learning theory provides generalization bounds for the learned quantizer.

#### 3.2 Capacity Analysis

The information capacity for each method:

$$C_{\text{CES}} = L_{\text{msg}} \text{ bits/token} \quad (9)$$

$$C_{\text{TGH}} = L \text{ bits/token} \quad (10)$$

$$C_{\text{KLQ}} = \log_2(C) \text{ bits/token} \quad (11)$$

## 4 Experimental Results

We evaluate on Mixtral-8x7B using PAWS [?] for adversarial paraphrase testing. Results show significant improvements over existing methods:

Our methods achieve 37-46% improvement in robustness (AUC) over existing approaches while maintaining comparable or better efficiency.

## 5 Conclusion

We introduced the first framework for MoE-native watermarking, representing a paradigm shift from treating MoE models as generic dense networks to exploiting their unique sparse computation patterns. Our three complementary methods

Method	AUC	$\Delta\text{PPL}$	ms/tok	bits/tok
Kirchenbauer et al.	0.65	0.8	2.1	1.0
RW-LSH (Patent)	0.58	0.9	2.3	8.0
<b>CES+ECC</b>	<b>0.89</b>	<b>0.7</b>	<b>2.2</b>	<b>6.0</b>
<b>TGH</b>	<b>0.92</b>	<b>1.1</b>	<b>2.8</b>	<b>8.0</b>
<b>KLQ</b>	<b>0.94</b>	<b>0.6</b>	<b>2.0</b>	<b>4.0</b>

Table 1: Performance comparison across key metrics. Our MoE-native methods show superior robustness while maintaining competitive efficiency.

demonstrate superior robustness against paraphrase attacks while maintaining competitive efficiency. This work opens new directions for architecture-aware security in sparse neural networks.

## Acknowledgments

We thank the anonymous reviewers for their valuable feedback.

## Availability

Code and datasets will be made publicly available upon publication.  
plain