

ZKPROV: A Zero-Knowledge Approach to Dataset Provenance for Large Language Models

Mina Namazi
Case Western Reserve University
Cleveland, Ohio, USA
mxn559@case.edu

Alexander Nemecek
Case Western Reserve University
Cleveland, Ohio, USA
ajn98@case.edu

Erman Ayday
Case Western Reserve University
Cleveland, Ohio, USA
exa208@case.edu

Abstract

As the deployment of large language models (LLMs) grows in sensitive domains, ensuring the integrity of their computational provenance becomes a critical challenge, particularly in regulated sectors such as healthcare, where strict requirements are applied in dataset usage. We introduce ZKPROV, a novel cryptographic framework that enables zero-knowledge proofs of LLM provenance. It allows users to verify that a model is trained on a reliable dataset without revealing sensitive information about it or its parameters. Unlike prior approaches that focus on complete verification of the training process (incurring significant computational cost) or depend on trusted execution environments, ZKPROV offers a distinct balance. Our method cryptographically binds a trained model to its authorized training dataset(s) through zero-knowledge proofs while avoiding proof of every training step. By leveraging dataset-signed metadata and compact model parameter commitments, ZKPROV provides sound and privacy-preserving assurances that the result of the LLM is derived from a model trained on the claimed authorized and relevant dataset. Experimental results demonstrate the efficiency and scalability of the ZKPROV in generating this proof and verifying it, achieving a practical solution for real-world deployments. We also provide formal security guarantees, proving that our approach preserves dataset confidentiality while ensuring trustworthy dataset provenance.

Keywords

Verifiable Large Language Model, Zero-Knowledge Proofs, Dataset Provenance

1 Introduction

Large language models (LLMs) are integrated into critical decision-making processes, including healthcare diagnostics [5], financial risk assessment [2], and legal services [3]. While LLMs are increasingly integrated in high-stakes applications, they introduce significant challenges in verifying their computational integrity [4, 20]. These verification challenges are particularly critical in regulated domains where trustworthiness is legally mandated.

To understand where verification efforts are focused, recent research categorizes machine learning integrity concerns into three core areas: (i) inference verification to validate model outputs, (ii) training process verification to confirm proper algorithm execution, and (iii) training data verification to ensure models are trained

on authorized datasets. Malicious actors may compromise these aspects to reduce computational costs or introduce privacy vulnerabilities, undermining trust in deployed systems, particularly when computations are outsourced due to privacy concerns or resource limitations [30].

Several cryptographic approaches have emerged to address different aspects of machine learning integrity. Secure Multi-party Computation (SMC) enables confidential, collaborative training but lacks scalability and requires constant connectivity among participants [16, 31]. Trusted Execution Environments (TEEs) offer hardware-based integrity guarantees but remain vulnerable to side-channel attacks and rely on centralized trust assumptions [29]. Homomorphic Encryption (HE) allows computation on encrypted data, but incurs prohibitive performance overhead for complex models [11].

Zero-knowledge proofs (ZKPs) demonstrate particular promise for verifiable machine learning. For example, systems such as VeriML [33], zkCNN [21], zkLLM [28], and zkGPT [22] focus on verifying the training process or inference correctness. However, these approaches do not address a critical concern regarding training data provenance. That is, the ability to cryptographically prove that a model was trained on a specific, authorized dataset without revealing the dataset itself.

This paper addresses a critical gap in current zero-knowledge-based approaches. While existing systems [7–9, 13, 15, 21–23, 28, 32] can prove the computational correctness of training or inference, they cannot cryptographically verify that a model was trained on a specific, authorized dataset without revealing the dataset itself or details of the model architecture. This capability is essential in settings where regulatory compliance demands that models draw only from relevant data sources, including credible information about the corresponding prompts. For example, if a medical professional queries an LLM with, “What is the life expectancy of patients diagnosed with breast cancer in this hospital between 1985 and 2025?”, the model’s response must be derived from the institution’s authorities clinical datasets containing the individual’s data with breast cancer in that hospital, not from public or external sources.

This work introduces a novel framework called ZKPROV that bridges this gap by providing dataset provenance verification without the computational burden of proving the entire training process. Our approach generates verifiable proofs that cryptographically bind a model response to its training dataset while maintaining the privacy of both the dataset and model parameters. Unlike existing solutions that either focus solely on inference correctness or require prohibitively expensive complete training verification, our framework offers a practical middle ground: statistical guarantees that the model could not have reached its state without exposure

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.
© 2025 Copyright held by the owner/author(s).
<https://doi.org/XXXXXX.XXXXXXX>



to the claimed training data, without proving every step of how it reached that state. Our main **contributions** are listed as follows.

- We design a privacy-preserving framework that cryptographically binds LLM responses to authorized training datasets without revealing dataset contents or model internals, achieving provenance verification with complete parameter confidentiality.
- We integrate efficient cryptographic schemes, including the ZK recursive proof system with hierarchical commitment structures, enabling sublinear scaling ($O(n \log n)$) for single-dataset scenarios and succinct verification of dataset inclusion across multiple transformer layers.
- We demonstrate practical scalability with sub-second proof generation for LLMs and establish performance thresholds for multi-dataset deployments. We achieve comparable efficiency to state-of-the-art zero-knowledge LLM verifications while maintaining privacy guarantees.
- We formally prove that the proposed framework is sound and preserves the privacy of the training dataset and model parameters in LLM, demonstrating robustness against known attacks.

The proposed approach represents a significant advancement in verifiable machine learning, specifically in LLMs. It offers a practical solution to the critical challenge of dataset provenance verification while preserving the privacy requirements essential for sensitive applications. By enabling zero-knowledge verification of dataset provenance, our framework empowers organizations to deploy LLMs in sensitive domains confidently, ensuring compliance standards while preserving privacy.

The remainder of this paper is organized as follows: Section 2 reviews related work in verifiable machine learning and zero-knowledge proofs. Section 3 presents the required background to design our framework in Section 4. Section 5 analyses the security and privacy of our proposed protocol. Section 6 presents our experimental evaluation. Section 7 discusses limitations and future work, and Section 8 concludes the paper.

2 Related Work

The ZKPROV occupies a unique position in verifiable machine learning by addressing dataset provenance rather than computational correctness. This section reviews related efforts in verifiable machine learning, with a focus on zero-knowledge proofs (ZKPs).

Traditional approaches like ZEN [8] focus on proving that inference computations are performed correctly on committed model parameters, achieving verification times in the range of seconds to minutes for large networks. ZEN provides efficient inference verification for large neural networks by introducing quantization-friendly encodings to reduce proof cost. Verification has also been extended to large language models (LLMs), such as in zkLLM [28], which optimizes the arithmetization of non-linear operations such as attention mechanisms and activation functions. This system for zkLLM introduces specialized components to handle transformer-specific operations at scale. These techniques allow a prover to demonstrate that a specific input-output pair was produced by a committed model, without exposing internal model details.

Additional efforts have extended inference verification to other model classes, including convolutional neural networks [21] and decision trees [32], using model-specific circuit optimizations.

Despite these advances, all of these systems share a key limitation where they verify inference correctness *after* training but offer no cryptographic assurance regarding *which dataset* was used during training. These approaches implicitly assume that all training data is authorized for use. This limitation is especially problematic in domains requiring dataset auditability and regulatory compliance, such as healthcare, where outputs must be derived exclusively from institutionally approved datasets.

Beyond verifying model outputs, recent work also explored verifying the training process using ZKPs. These systems introduce protocols that prove a model was trained correctly, typically via gradient descent, on a committed dataset [1, 10]. Such approaches offer verifiability over the whole computational training process by requiring proofs for each optimization step. However, as model sizes and training complexity grow, these methods become increasingly computationally expensive and are generally unscalable for large models such as LLMs.

Moreover, in these systems, no explicit cryptographic guarantees are provided regarding the authenticity of training data from approved sources. Even if the training procedure is proven to follow a prescribed algorithm, there is no mechanism to bind the resulting model to a specific, authorized dataset.

We introduce a new category in the landscape of verifiable machine learning with privacy-preserving dataset provenance for LLMs. Unlike prior work, our framework does not require complete training process verification. Instead, it provides ZKPs that a model was trained or fine-tuned using a specific, authenticated dataset, while keeping its contents confidential.

Contrasting our approach with retrieval-augmented generation (RAG), which retrieves relevant information from a hosted dataset at inference time, is essential. While the dataset used in RAG may be approved for access, RAG systems do not guarantee that the underlying LLM was trained solely on authorized data. In regulated settings, it is often necessary to ensure that users interact not just with approved retrieval sources but with models whose behavior is fully constrained by authorized training datasets. Our framework addresses this need by enabling cryptographic verification that a model response originated from a model bound to such datasets.

Although RAG techniques could potentially be integrated into the ZKPROV framework, our current focus is verifying dataset provenance for trained models, independent of how the training was performed. While our evaluation fine-tunes a model to demonstrate domain relevance, the ZKPROV protocol determines whether the model is pre-trained, fine-tuned, or trained from scratch. This choice reflects a practical deployment model where training and serving a model with verified dataset provenance is simpler than securely hosting and governing access to the datasets themselves at inference time. We further discuss this integration opportunity in Section 7.

To the best of our knowledge, ZKPROV is the first system to bridge the gap between data regulation and cryptographic LLM auditability. It enables stakeholders to verify that a model response could only have been generated by a model trained on approved

data, without revealing information about the model, dataset, or training process.

3 Background

We deploy state-of-the-art cryptographic schemes that efficiently prove computations' correctness and address the challenge of securely verifying LLM provenance.

We briefly introduce the basic building blocks and encryption schemes in our proposed privacy-preserving LLM provenance verification framework.

3.1 Large Language Model Primitives

Large Language Models (LLMs) are the main components in our proposed framework. An LLM is treated as a parameterized function $LLM(W, p)$ that maps a natural language prompt p to a response r , where W denotes the model's weights, which may evolve through fine-tuning. Given the sensitivity of downstream applications, such as healthcare, we must ensure that the final model response r is provably linked to a specific, authorized dataset.

Fine-Tuning. We define the following function to refer to the fine-tuning process for the initial LLM model weights W_0 and the dataset: $W \leftarrow \text{FINETUNE}(D, W_0, H)$, Where $H = (\eta, B, E, O)$ is specified as the learning rate, batches, epochs, and the optimizer, respectively, to define the acceptable training configurations for regulatory compliance.

Privacy-Preserving Dataset Retrieval. We define a retrieval function that selects the optimal dataset given a prompt and its associated attributes to ensure that the model is derived from an authorized and contextually relevant dataset. Let p be a natural language query, and Att_p denote its associated attribute set (e.g., clinical topic, patient demographics). The retrieval function is defined as:

$$D_{i^*} \leftarrow \text{MATCH}(p, Att_p, \{(D_i, Att_i)\}_{i=1}^m)$$

where D_{i^*} is the selected dataset such that MATCH maximizes semantic and attribute relevance between the prompt attributes Att_p and the dataset attributes Att_i . The selection mechanism may compute similarity based on keyword overlap, ontology-aware mappings¹ (e.g., shared UMLS codes), or embedding-based distance², depending on deployment constraints. This ensures that downstream model responses are provably linked to the most appropriate dataset in a privacy-preserving manner.

If no dataset exceeds a predefined similarity threshold, the system outputs a null result, indicating that no suitable dataset is available. This prevents responses from being generated using irrelevant or unauthorized data, thereby preserving both privacy and provenance integrity.

3.2 Cryptographic Primitives

We describe the core cryptographic schemes and their hardness assumptions, on which we rely when developing the proposed

privacy-preserving LLM provenance verification framework. Our scheme is built upon elliptic curves, as described in the following.

Let \mathbb{G}_1 and \mathbb{G}_2 be elliptic curve groups of prime order q , and let \mathbb{G}_T denote the target group. A bilinear map $e : \mathbb{G}_1 \times \mathbb{G}_2 \rightarrow \mathbb{G}_T$ satisfies:

- Bilinearity: For all $g_1 \in \mathbb{G}_1$, $g_2 \in \mathbb{G}_2$, and $a, b \in \mathbb{Z}_q$, the following equation holds, $e(g_1^a, g_2^b) = e(g_1, g_2)^{ab}$.
- Non-degeneracy: $e(g_1, g_2) \neq 1$, where g_1 and g_2 are generators of \mathbb{G}_1 and \mathbb{G}_2 , respectively.

3.2.1 Zero-Knowledge Succinct Non-Interactive Argument of Knowledge (zk-SNARK). A zk-SNARK is a cryptographic proof system that allows a prover to convince a verifier that a statement $x \in \mathcal{L}_R$ is valid with respect to a relation R , without revealing any auxiliary information (i.e., the witness ω).

A zk-SNARK operates on a computation represented as a constraint system, which expresses the computation as a set of polynomial constraints on a computation represented as a *Rank-1 Constraint System (R1CS)*, which is a way to express a computation as a set of quadratic constraints. It consists of three main components as follows. For our LLM provenance verification framework, we deploy HyperNova [18], an efficient recursive proof system, allowing us to express complex LLM computations involving high-degree polynomial operations. HyperNova consists of $\text{HN} = \{\text{ZK.SETUP}, \text{ZK.PROVE}, \text{ZK.VERIFY}, \text{ZK.FOLD}\}$, described as follows.

- $\text{ZK.SETUP}(1^\lambda, R) \rightarrow pp_{zk}$: On a given security parameter λ and a relation R , generate common reference strings (*crs*) for the underlying polynomial commitment scheme, encoding the constraint matrices $\tilde{M}_1, \dots, \tilde{M}_l$ as sparse multilinear polynomials, and outputs $pp_{zk} = \{crs, pk, vk\}$.
- $\text{ZK.PROVE}(pk, x, \omega) \rightarrow \pi$: Input the proving key pk , a public statement x , and a private witness ω , and generate a proof π by executing a multi-folding scheme, which reduces multiple proof instances into a single one.
- $\text{ZK.VERIFY}(vk, x, \pi) \rightarrow 0/1$: Take the verification key vk , the public input x , and the proof π to checks the validity of the proof according to following of bilinear pairings equation: $e(C, g) = e(\pi, g^{crs})$. Where $C = g^{f(\alpha)}$ is the commitment to the polynomial $f(X)$.
- $\text{ZK.FOLD}(pk, (U_1, w_1), (U_2, w_2)) \rightarrow (U', w')$: This folding algorithm combines two instances (U_1, w_1) and (U_2, w_2) into a single one (U', w') . It compresses multiple computational steps using a random linear combination:

$$U' = U_1 + r \cdot U_2, \quad w' = w_1 + r \cdot w_2,$$

Where $r \in \mathbb{F}_p$ is a random challenge the verifier generates. The folding process ensures that the resulting instance (U', w') is satisfiable if and only if both original instances (U_1, w_1) and (U_2, w_2) are satisfiable. This enables incremental proof generation for multi-layer computations in neural network architectures.

Unlike traditional SNARKs that require trusted setup ceremonies, HyperNova uses universal setup parameters that can be reused across different constraint systems, making it practical for evolving LLM architectures. The recursive composition ensures that the proof size remains constant, regardless of the number of steps, significantly improving scalability. The verifier's work is logarithmic

¹For example, concept alignment using clinical ontologies such as the Unified Medical Language System (UMLS).

²For example, computing cosine similarity between vectorized representations of prompts and dataset metadata using pre-trained language models.

in the size of the constraint system, $O(\log m)$, enabling efficient verification for large-scale LLM computations. This property is ideal for applications involving iterative or modular computations, such as verifiable machine learning systems.

In the following, we define the security properties of the zk-SNARK scheme.

Negligible Function. A function $\nu : \mathbb{N} \rightarrow \mathbb{R}_{\geq 0}$ is negligible if for every positive polynomial $p(\cdot)$, there exists an N_p such that for all $n > N_p$, $\nu(n) < \frac{1}{p(n)}$.

Soundness. A zk-SNARK proof system satisfies computational soundness if, for every probabilistic polynomial-time (PPT) adversary Adv :

$$\Pr [\text{ZK.VERIFY}(vk, x, \pi) = 1 \wedge x \notin \mathcal{L}_R] \leq \nu(\lambda),$$

where the probability is over the randomness of ZK.SETUP , Adv , and ZK.VERIFY .

Zero-Knowledge. A zk-SNARK proof system satisfies computational zero-knowledge if there exists a PPT simulator \mathcal{S} such that for every PPT verifier \mathcal{V}^* and every $x \in \mathcal{L}_R$, the distributions of the proofs $\pi_1 = \text{ZK.PROVE}(pk, x, w_1)$ and $\pi_2 = \text{ZK.PROVE}(pk, x, w_2)$ are identical.

$$\Pr[Adv(\pi_1) = 1] - \Pr[Adv(\pi_2) = 1] \leq \nu(\lambda),$$

where $\pi_1 = \text{ZK.PROVE}(pk, x, w_1)$ and $\pi_2 = \text{ZK.PROVE}(pk, x, w_2)$, and the probabilities are over the randomness of ZK.PROVE and the adversary Adv .

Privacy. A zk-SNARK proof system satisfies privacy, if there exists a PPT simulator \mathcal{S} such that for every PPT verifier \mathcal{V}^* and every $x \in \mathcal{L}_R$, the verifier's view during the interaction with the prover is computationally indistinguishable from the simulator's output. Formally:

$$\text{View}_{\mathcal{V}^*}(pk, x) \approx_c \mathcal{S}(pk, x).$$

Two ensembles of probability distributions $\{X_\lambda\}_{\lambda \in \mathbb{N}}$ and $\{Y_\lambda\}_{\lambda \in \mathbb{N}}$ are computationally indistinguishable, denoted $X_\lambda \approx_c Y_\lambda$, if for every PPT algorithm \mathcal{D} (the distinguisher), there exists a negligible function $\nu(\cdot)$ such that for all sufficiently large $\lambda \in \mathbb{N}$:

$$\left| \Pr[\mathcal{D}(1^\lambda, x) = 1 : x \leftarrow X_\lambda] - \Pr[\mathcal{D}(1^\lambda, y) = 1 : y \leftarrow Y_\lambda] \right| \leq \nu(\lambda).$$

3.2.2 Commitment Scheme. A commitment scheme is a cryptographic primitive that allows a party to commit to a value while keeping it hidden. They are usually hiding; no information is revealed about the committed value. They are also binding, meaning that once a commitment is created, it is computationally infeasible to change it to a different value. The message and opening information can later be revealed.

We deploy the Kate-Zaverucha-Goldberg (KZG) [17], a commitment scheme designed for polynomials. It allows committing to a polynomial $f(X)$. It later proves the evaluation of the polynomial at any point z without revealing the entire polynomial. It consists of $\text{KZG} = \{\text{C.SETUP}, \text{C.COMMIT}, \text{C.OPEN}, \text{C.VERIFY}\}$, described as follows.

- **C.SETUP** $(1^\lambda, d) \rightarrow pp$: Sample $\alpha \leftarrow \mathbb{F}_p$, generate public parameter $pp = (g, g^\alpha, g^{\alpha^2}, \dots, g^{\alpha^d}) \in \mathbb{G}_1^{d+1}$.

- **C.COMMIT** $(pp, f(X)) \rightarrow C$: For polynomial $f(X)$, generate commitment $C = \prod_i (g^{\alpha^i})^{f_i}$. This commitment is a single group element, making it compact regardless of the polynomial's degree.
- **C.OPEN** $(pp, f(X), z) \rightarrow \pi$: Generate a proof π for evaluation $f(z)$ at point z :

$$\pi = g^{q(\alpha)} \quad \text{where } q(X) = \frac{f(X) - f(z)}{X - z}.$$

The proof π ensures that the prover knows the entire polynomial $f(X)$ without revealing it. The quotient polynomial $q(X)$ guarantees the correctness of $f(z)$.

- **C.VERIFY** $(pp, C, z, y, \pi) \rightarrow 0/1$: Check correctness of $f(z)$ using pairing:

$$e(C_f / g^{f(z)}, g) = e(\pi, g^{\alpha-z}).$$

This pairing equation ensures that the committed polynomial evaluates correctly at z , and the verifier does not need to know α or $f(X)$, preserving privacy and efficiency.

The KZG scheme is homomorphic, enabling proof aggregation across various constraint systems. This allows for verification of a single evaluation independent of the polynomial's size or degree, significantly reducing computational and verification overhead for complex systems like LLMs.

3.2.3 Tree-Based Commitment Scheme. We deploy Reckle Trees [26] to organize the elements of the datasets due to their efficiency in supporting batch updates. Particularly, in healthcare scenarios, where datasets frequently require new patient records, updated treatment protocols, or regulatory compliance updates, proving dataset membership for multiple data points simultaneously becomes highly efficient.

3.2.4 Boneh-Lynn-Shacham (BLS). We use Boneh-Lynn-Shacham (BLS) [6] signature scheme, where a signature is required on the data. The BLS provides short signatures with efficient verification over elliptic curve groups with pairing operations. The scheme offers strong unforgeability under chosen message attacks. It enables signature aggregation for batch verification scenarios, which are critical for minimizing communication overhead and verifying multiple datasets signed by various authorities simultaneously. The signature generation and verification are below.

- **BLS.SIGN** $(sk, m) \rightarrow \sigma$: Takes private key sk and message m , computes the signature σ .
- **BLS.VERIFY** $(pk, m, \sigma) \rightarrow \{0, 1\}$: Takes public key pk , message m , and signature σ , verifies it, and outputs 1 if the verification succeeds, 0 otherwise.

4 Proposed Scheme

The proposed ZKPROV framework establishes a cryptographically robust method for verifying that a large language model (LLM) is trained on specific, authorized datasets. While our framework applies generally to any LLM training process, we focus on fine-tuning in our implementation since domain-specific applications require adapting pre-trained models to particular datasets. Without loss of generality, we use fine-tuning terminology, though the underlying cryptographic mechanisms apply equally to training from

scratch or any other model adaptation process. A core objective is to achieve this verification without compromising the privacy of the dataset contents or the model’s parameters. This section represents the proposed scheme’s setting, threat model, overview, and details. We provide an overview of the proposed scheme in Figure 1.

4.1 System Model and Settings

The proposed protocol comprises an authority \mathcal{CA} (e.g., health-care institution) that owns and authenticates datasets; a prover \mathcal{P} , the LLM service provider, that receives user queries, generates responses, and corresponding proofs; and a user \mathcal{U} who submits queries and can share responses with verifiers to check the authenticity of the underlying training dataset.

The Authority \mathcal{CA} owns and maintains a collection of datasets $D = \{D_1, \dots, D_m\}$. Each D_i is characterized by attribute sets $Att_i = \{att_1, \dots, att_k\}$, reflecting demographic, domain-specific, and contextual information (e.g., $Att_1 = \{\text{“cancer”}, \text{“female”}, \text{“age_65+”}\}$). It also includes administrative identities as $Id_i = \{id_1, \dots, id_k\}$.

A prompt p from user \mathcal{U} consists of natural language text pertaining to a downstream task. A response r consists of a natural language text answering the query, accompanied by a cryptographic proof π that verifiably links the response to the authenticated D_i through zero-knowledge protocols, ensuring response authenticity without compromising model privacy or dataset confidentiality.

4.2 Threat Model

The authority \mathcal{CA} is assumed to be honest in authenticating genuine datasets and generating public parameters for the entire framework. The user \mathcal{U} is semi-honest, i.e., they follow the protocol’s steps but might be curious to gain unauthorized information from the interaction’s outputs. The LLM service provider \mathcal{P} can be malicious and may not follow the protocol’s instructions, and might attempt to use unauthorized data or forge the outputs of the protocols.

A malicious \mathcal{P} may attempt to train models on unauthorized or modified datasets while claiming to use authenticated ones. They may attempt to use datasets not signed by \mathcal{CA} or imply that they have been tampered with after authentication. Also, it may attempt to commit to incorrect weight differences to reflect the model’s behavior inaccurately. Finally, the provider may try to forge the proofs to break the cryptographic link between responses and authenticated datasets. Malicious users cannot inject false data, but might be curious to obtain information about the datasets’ content while verifying the proofs.

Formally, in the ZKPROV framework, the following security properties are guaranteed.

- **Zero-Knowledge (ZK):** The proofs of provenance must not reveal any information about the confidential model weight differences or the underlying training dataset, beyond what is explicitly disclosed through the commitments to dataset metadata.
- **Soundness:** The ZKPROV system must be sound, meaning that a computationally bounded malicious prover cannot convince an honest verifier of a false statement regarding the dataset provenance or the model’s computational integrity, except with negligible probability.

- **Dataset Exposure Binding (DEB):** It ensures that a valid proof, attesting to using a specific authorized dataset, implies that the model’s state was derived under the influence of these. This property prevents a malicious prover from successfully claiming provenance from an authorized dataset while using an unauthorized one.

4.3 Overview

The ZKPROV framework establishes a verifiable link between LLM responses and authenticated datasets through an efficient zero-knowledge proof system. The authority \mathcal{CA} owns multiple datasets and cryptographically authenticates them using metadata and digital signatures. These datasets are stored in cryptographic structures, ensuring both integrity and authorization.

The provider \mathcal{P} uses these authenticated datasets to fine-tune a base model and commits to all sensitive protocol components, including the model weights and dataset-specific parameters. By leveraging ZK proofs, the provider ensures that sensitive information, such as dataset contents and model parameters, remains confidential while maintaining verifiable integrity.

At query time, when a user \mathcal{U} submits a prompt, the provider selects the most relevant dataset, generates a response using the model, and produces a proof demonstrating that the response was derived from a model trained on authenticated data. The proof links the response, the prompt, and the dataset while preserving the privacy of the corresponding dataset.

The verification phase ensures that all cryptographic relationships in the proof are valid. The user’s first check is to ensure that the dataset metadata carries a valid signature from the authority \mathcal{CA} , proving that the metadata is authorized and originates from a trusted source. Next, the user verifies that the provider correctly computed the binding values, which link the model’s weights to challenge vectors derived from the authenticated dataset metadata. The user also validates that the weight differences between the base and updated models are consistent with the transformations induced by the authenticated dataset. Finally, the verification process ensures that the proof corresponds to the specific query-response interaction. The query transcript binds the proof to the user’s prompt and response, preventing reuse of the proof in other contexts.

By combining privacy-preserving zero-knowledge proofs with large language model primitives, ZKPROV enables robust provenance verification for LLM responses.

The proposed privacy-preserving, verifiable LLM framework ZKPROV comprises 3 main algorithms that work together to establish and verify the required verification chain.

ZKPROV = {SETUP, PROVE, VERIFY} as defined in the following.

- **SETUP**($1^\lambda, W_0$) $\rightarrow (C, \Omega)$: Given security parameter λ and base model W_0 , generates the necessary cryptographic parameters for all participants, performs dataset authentication and model fine-tuning, and produces public commitments C and private witnesses Ω for all sensitive protocol components.
- **PROVE**(C, Ω, p) $\rightarrow (r, \pi)$: Given commitments C , witnesses Ω , and user prompt p , selects the optimal dataset, generates response r using the corresponding model, and produces a proof π demonstrating that the response derives from a model trained on authenticated data.

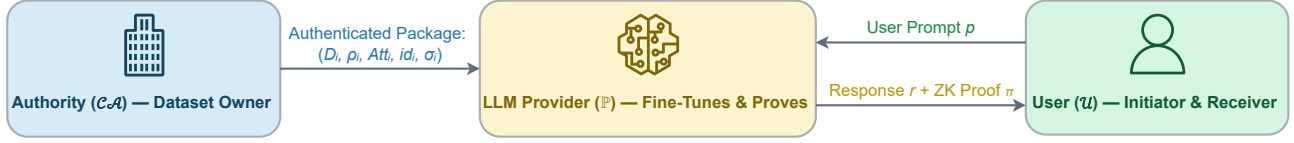


Figure 1: High-level ZKPROV framework’s flow. The figure illustrates the interaction between the Authority (\mathcal{CA}), the LLM Provider (\mathcal{P}), and the User (\mathcal{U}). The Authority certifies dataset authenticity and transmits signed metadata to the provider, who fine-tunes a base model and produces provenance-aware responses. The User submits a prompt and receives a response along with a zero-knowledge proof attesting to the dataset’s validity. The labeled arrows highlight the authenticated dataset handoff, prompt input, and proof-backed output.

- $\text{VERIFY}(p, r, \pi, C) \rightarrow \{\text{ACCEPT}, \text{REJECT}\}$: Given prompt p , response r , proof π , and commitments C , validates all cryptographic relationships to confirm that the response was derived from a model trained on the authenticated datasets without learning sensitive information about dataset contents or model parameters.

4.4 Detailed Description of ZKPROV

This section provides proposed ZKPROV scheme’s details, comprising setup and preprocessing, proof construction, and verification phases.

4.4.1 Setup and Preprocessing. The authority \mathcal{CA} generates keys and public parameters of the framework and distributes them among the parties to be the inputs of all the components. It generates the common reference string (crs) that contains domain separators κ_1, κ_2 , and κ_3 that ensure cryptographic separation between different protocol components, preventing collision attacks across operational domains.

\mathcal{CA} stores the datasets in a Reckle Tree T_i (Sec. 3.2.3) and their roots ρ_i as commitments to the entire dataset. For each dataset, the authority constructs metadata $m_i = (\rho_i, \text{Atti}, \text{idi})$ that combines the ρ_i with semantic attributes and administrative identifiers. The authority then generates BLS signatures σ_i (Sec. 3.2.4), on each metadata package m_i using their secret key $sk_{\mathcal{CA}}$, creating unforgeable authentication tokens that cryptographically bind each dataset to its authorized usage constraints. They transmit the authentication package $AP = (pp, D_i, m_i, \sigma_i, T_i)$ to the \mathcal{P} .

\mathcal{P} receives AP from the authority and verifies the signature using the \mathcal{CA} ’s public key to ensure dataset authenticity. If the verification is successful, the \mathcal{P} fine-tunes the base model W_0 using their specified hyperparameters $H = (\eta, B, E, O)$ to obtain dataset-specific weights $W_i \leftarrow \text{FINETUNE}(D_i, W_0, H)$. The provider then computes layer-wise weight differences $\Delta W_{i,j} = W_{i,j} - W_{0,j}$ for each layer $j \in \{1, \dots, N\}$, capturing the specific transformations on dataset D_i .

For each dataset, the provider generates commitments to the authorities signatures and metadata by calling $C_{\sigma,i} \leftarrow \text{C.COMMIT}(\sigma_i, \omega_{\sigma,i})$ and $C_{m_i} \leftarrow \text{C.COMMIT}(m_i, \omega_{m_i})$ for the signature and metadata respectively, along with commitments to the weight differences $C_{\Delta W,i} \leftarrow \text{C.COMMIT}(\Delta W_i, \omega_{\Delta W,i})$ and the base model $C_{W_0} \leftarrow \text{C.COMMIT}(W_0, \omega_{W_0})$. The randomness values $\omega_{\sigma,i}$ and ω_{m_i} serve a critical role beyond standard commitment hiding, as they are used to generate the seed value $\text{seed}_i = \text{PRF}(\omega_{\sigma,i} \parallel \omega_{m_i} \parallel \kappa_2)$ that ensures the following challenge vectors can only be computed by solely possessing valid commitments to authenticated m_i ’s.

For each dataset i and layer j , the provider computes challenge vectors $v_{i,j} = \text{HASH}(\text{seed}_i \parallel j \parallel \kappa_1)$, where the seed incorporates the commitment randomness and the domain separator κ_2 ensures contextual uniqueness across different datasets. The $v_{i,j}$ generation mechanism maintains the non-interactive nature of the zero-knowledge proof system while ensuring unpredictability properties essential for security. This construction guarantees that $v_{i,j}$ are deterministically derivable during proof generation and verification, yet remain computationally unpredictable to any \mathcal{P} who does not possess the correct signature on authenticated m_i and the corresponding ω_{m_i} . The layer-specific domain separator κ_1 ensures that different layers produce distinct challenge vectors even for the same dataset, preventing cross-layer collision attacks.

The binding computation creates cryptographic fingerprints that uniquely identify the weight differences through inner product evaluations. For each dataset i and layer j , the provider calculates

$$B_{i,j} \leftarrow \langle \Delta W_{i,j}, v_{i,j} \rangle = \sum_{k=1}^d \Delta W_{i,j}[k] \cdot v_{i,j}[k].$$

Where d represents the dimensionality of the weight vectors in layer j . These binding values are cryptographically sound under the Schwartz-Zippel lemma [27], which guarantees that distinct multivariate polynomials evaluated at random points yield different results with overwhelming probability. The provider then commits to these binding values as $C_{B,i} \leftarrow \text{C.COMMIT}(\{B_{i,j}\}_{j=1}^N, \omega_{B,i})$.

The final commitment vector hides the actual values $m_i, \sigma_i, \Delta W_{i,j}, B_{i,j}$, and W_0 while enabling their use within the zero-knowledge proofs generation mechanism. $C = (C_{m_i}, C_{\sigma,i}, C_{\Delta W,i}, C_{B,i}, C_{W_0})$

The challenge vectors $v_{i,j}$ are not explicitly committed to as part of this public set C , as their correct derivation and usage is proven internally during the ZK proof based on the witnesses that open C_{m_i} and $C_{\sigma,i}$. This design choice reduces the public commitment size while maintaining the ability to verify that challenge vectors were computed correctly using authenticated inputs. The provider stores all commitment opening information $\Omega = (\omega_{m_i}, \omega_{\sigma,i}, \omega_{\Delta W,i}, \omega_{B,i}, \omega_{W_0})$ as private witnesses to be used during proof generation to demonstrate knowledge of the committed values and their relationships as shown in Alg. 1.

4.4.2 Proof Construction. The user \mathcal{U} submits a prompt p , triggering the embedded retrieval algorithm that selects the optimal dataset. The model provider identifies the most relevant dataset by calling $D_{i^*} \leftarrow \text{MATCH}(p, \text{Atti}_p, \{(D_i, \text{Atti}_i)\}_{i=1}^m)$.

where the matching function evaluates the prompt attributes Atti_p against the dataset attributes $\{\text{Atti}_i\}_{i=1}^m$ to identify the most contextually relevant dataset. This selection process ensures that

the response is generated using the model weights W_i^* that were specifically trained on the dataset most appropriate for the query domain. For simplicity and avoiding complex notation, we refer to the selected dataset index as i (i.e., $i \leftarrow i^*$) for the remainder of the proof construction process.

Algorithm 1 Setup and Preprocessing Phase of ZKPROV

```

1: procedure SETUP( $\{D_i\}_{i=1}^m, W_0, \text{crs}$ )
    The  $\mathcal{CA}$  generates  $pp$ , and authenticates datasets:
2:   for  $i = 1$  to  $m$  do
3:     Generates  $AP = (D_i, m_i, \sigma_i, T_i)$  to  $\mathcal{P}$ 
4:   end for
    The  $\mathcal{P}$  verifies authentication and fine-tunes the model:
5:   for  $i = 1$  to  $m$  do
6:     if Signature  $\sigma_i$  on  $m_i$  is valid then
7:        $W_i \leftarrow \text{FINETUNE}(D_i, W_0, H)$ 
8:       for  $j = 1$  to  $N$  do
9:          $\Delta W_{i,j} \leftarrow W_{i,j} - W_{0,j}$ 
10:      end for
11:     end if
12:   end for
    The  $\mathcal{P}$  generates commitments:
13:    $\omega_{W_0} \xleftarrow{\$} \mathbb{F}_p$ 
14:    $C_{W_0} \leftarrow \text{C.COMMIT}(W_0, \omega_{W_0})$ 
15:   for  $i = 1$  to  $m$  do
16:      $\omega_{\sigma,i}, \omega_{m,i}, \omega_{\Delta W,i}, \omega_{B,i} \xleftarrow{\$} \mathbb{F}_p$ 
17:      $C_{\sigma,i} \leftarrow \text{C.COMMIT}(\sigma_i, \omega_{\sigma,i})$ 
18:      $C_{m,i} \leftarrow \text{C.COMMIT}(m_i, \omega_{m,i})$ 
19:      $C_{\Delta W,i} \leftarrow \text{C.COMMIT}(\Delta W_i, \omega_{\Delta W,i})$ 
20:   end for
    The  $\mathcal{P}$  generates challenge vectors and binding values:
21:   for  $i = 1$  to  $m$  do
22:      $\text{seed}_i \leftarrow \text{PRF}(\omega_{\sigma,i} || \omega_{m,i} || \kappa_2)$ 
23:     for  $j = 1$  to  $N$  do
24:        $v_{i,j} \leftarrow \text{HASH}(\text{seed}_i || j || \kappa_1)$ 
25:        $B_{i,j} \leftarrow \langle \Delta W_{i,j}, v_{i,j} \rangle$ 
26:     end for
27:      $C_{B,i} \leftarrow \text{C.COMMIT}(\{B_{i,j}\}_{j=1}^N, \omega_{B,i})$ 
28:   end for
    The  $\mathcal{CA}$  assembles the final commitments and witnesses:
29:    $C \leftarrow \{C_{m,i}, C_{\sigma,i}, C_{\Delta W,i}, C_{B,i}\}_{i=1}^m \cup \{C_{W_0}\}$ 
30:    $\Omega \leftarrow \{\omega_{m,i}, \omega_{\sigma,i}, \omega_{\Delta W,i}, \omega_{B,i}\}_{i=1}^m \cup \{\omega_{W_0}\}$ 
31:   Return  $(C, \Omega)$ 
32: end procedure
    
```

The \mathcal{P} utilizes the model weights corresponding to the selected dataset and produces the LLM output $r = \text{LLM}(W_i, p)$

Using the dataset-specific model weights, ensuring that the response reflects the knowledge patterns embedded during fine-tuning on the authenticated dataset D_i . Since the response r is known to the verifier after transmission, it can be directly incorporated into the cryptographic transcript that binds the entire query-response interaction to the specific dataset used for generation. The provider creates a query transcript $\tau \leftarrow \text{HASH}(C_{m_i} || p || r || \kappa_3)$.

where κ_3 is the query binding separator defined in the crs . This transcript serves as a public digest that succinctly binds the committed dataset metadata, the user p , and the generated r into a single

cryptographic identifier that can be verified without revealing the underlying sensitive information.

The ZK proof construction requires demonstrating knowledge of valid openings for all commitments while proving the correctness of multiple relationships. The provider must prove that the committed metadata m_i carries an authentic signature from authority \mathcal{CA} , that the challenge vectors $v_{i,j}$ were correctly derived from the authenticated metadata and commitment randomness, and that the binding formula $B_{i,j}$ is computed correctly for each layer j . Additionally, the proof must establish that the weight differences $\Delta W_{i,j}$ are consistent with the model W_i derived from the base model W_0 and dataset D_i , and that the entire process aligns with the public query transcript τ .

The signature authenticity proof demonstrates that the committed signature σ_i is valid for the committed metadata m_i under the authority's public key. The provider constructs a zero-knowledge proof.

$$\pi_\sigma \leftarrow \text{ZK.PROVE}(pp_{zk}, (C_{m_i}, C_{\sigma,i}, pk_{\mathcal{CA}}), (\sigma_i, m_i, \omega_{m_i}, \omega_{\sigma,i}))$$

The proof ensures that only responses derived from models trained on authenticated datasets can be validated, preventing the use of unauthorized or tampered training data.

The binding correctness proof demonstrates that each binding value $B_{i,j}$ was computed as the correct inner product between the committed weight differences and the challenge vectors derived from the authenticated metadata.

The provider leverages HyperNova's recursive folding mechanism to aggregate all layer-wise binding constraints into a single succinct proof. The folding process begins with an empty state and iteratively incorporates each layer's binding constraint through

$$S_{B,j} \leftarrow \text{ZK.FOLD}(S_{B,j-1}, R_{B,i,j}, \omega_{B,j}).$$

Where $R_{B,i,j}$ represents the complete relation including both the inner product computation and the challenge vector derivation, and $\omega_{B,j} = (\Delta W_{i,j}, B_{i,j}, \omega_{\Delta W,i}, \omega_{B,i}, \omega_{m,i}, \omega_{\sigma,i})$ contains all necessary witness values. After processing all N layers, the provider generates the final recursive proof that succinctly demonstrates the correctness of all binding computations and challenge vector derivations with constant verification complexity.

$$\pi_B^{rec} \leftarrow \text{ZK.PROVE}(S_{B,N}).$$

The proof of weight consistency establishes that the committed weight differences correctly represent the transformation from the base model to the model for each layer. Since we have layer-wise weight differences $\Delta W_{i,j}$ for $j \in \{1, \dots, N\}$, the provider must use recursive folding to prove the relation $\Delta W_{i,j} = W_{i,j} - W_{0,j}$

Holds for all layers simultaneously.

The recursive weight difference proof begins with an empty state and iteratively incorporates each layer's weight consistency constraint through

$$S_{\Delta,j} \leftarrow \text{ZK.FOLD}(S_{\Delta,j-1}, R_{\Delta,i,j}, \omega_{\Delta,j}),$$

where $R_{\Delta,i,j}$ represents the relation $\Delta W_{i,j} = W_{i,j} - W_{0,j}$ and $\omega_{\Delta,j} = (W_{0,j}, W_{i,j}, \Delta W_{i,j}, \omega_{W_0}, \omega_{W,i}, \omega_{\Delta W,i})$ contains the layer-specific witness values and opening randomness.

After processing all N layers, the provider generates the final recursive weight consistency proof

$$\pi_{\Delta}^{rec} \leftarrow \text{ZK.PROVE}(S_{\Delta,N}),$$

That prevents malicious providers from claiming that arbitrary weight differences correspond to legitimate fine-tuning processes.

The query transcript consistency proof links the entire provenance chain to the specific query-response interaction, demonstrating that the committed dataset metadata used in the proof corresponds to the dataset selection for the given prompt and response. The provider shows that the query transcript $\tau = \text{HASH}(C_{m_i} \parallel p \parallel r \parallel \kappa_3)$ was correctly computed using the committed metadata for the selected dataset, the user's prompt, and the generated response. This proof

$$\pi_{\tau} \leftarrow \text{ZK.PROVE}(pp_{zk}, (C_{m_i}, p, r, \tau), (m_i, \omega_{m_i}))$$

ensures that the provenance verification is specific to the actual query-response pair rather than a generic proof that could be reused across different interactions.

The final proof aggregation combines all individual proof components into a single succinct zero-knowledge proof that can be efficiently verified. The provider assembles $\pi \leftarrow (\pi_{\sigma}, \pi_B^{rec}, \pi_{\Delta}, \pi_{\tau})$, containing all necessary proof elements to validate the complete provenance chain. This aggregated proof demonstrates that the provider knows secret witnesses that correctly open all public commitments and satisfy all underlying cryptographic relationships, while maintaining zero-knowledge properties that prevent the verifier from learning any information about the dataset contents, model parameters, or internal computational details beyond what is necessary for provenance verification. We provide a summary of the interactions in the Alg. 2

4.4.3 Verification phase. The user \mathcal{U} receives the response r and the aggregated proof π from the provider, along with the public commitments C , the query transcript τ , and validates each component of the provenance chain to ensure that all cryptographic relationships hold correctly and that the response genuinely originates from a model trained on authority-authenticated data.

The signature authenticity verification confirms that the committed signature σ_i is valid for the committed metadata m_i under the authority's public key pk_{CA} . The verifier executes

$$\{0, 1\} \leftarrow \text{ZK.VERIFY}(pp_{zk}, (C_{m_i}, C_{\sigma,i}, pk_{CA}), \pi_{\sigma})$$

to ensure that the dataset used for training was properly authenticated by the designated authority rather than being an unauthorized or tampered dataset.

Then the \mathcal{U} validates the recursive binding among all the layers, ensuring their correct computation using challenge vectors derived from authenticated metadata. The verifier calls

$$\{0, 1\} \leftarrow \text{ZK.VERIFY}(pp_{zk}, S_{B,N}, \pi_B^{rec}),$$

where $S_{B,N}$ represents the final folded state containing all N binding constraints. This verification confirms that each binding value is computed as the correct inner product of $\Delta W_{i,j}$ and $v_{i,j}$. Since this proof is generated recursively using the folding mechanism of the HyperNova proof system, its verification is constant-time regardless of the number of model layers, making the system scalable to large language models with hundreds of layers.

Algorithm 2 Proof Construction Phase of ZKPROV

```

1: procedure PROVE( $C, \Omega, p$ )
2:   Input: Commitments  $C$ , Witnesses  $\Omega$ , Prompt  $p$ 
   The  $\mathcal{P}$  constructs the proofs:
3:    $i^* \leftarrow \text{MATCH}(p, \text{Att}_p, \{(D_i, \text{Att}_i)\}_{i=1}^m)$ 
4:    $i \leftarrow i^*$ 
5:    $r \leftarrow \text{LLM}(W_i, p)$ 
6:    $\tau \leftarrow \text{HASH}(C_{m_i} \parallel p \parallel r \parallel \kappa_3)$ 
7:    $\text{witness}_{\sigma} \leftarrow (\sigma_i, m_i, \omega_{m_i}, \omega_{\sigma,i})$ 
8:    $\text{public}_{\sigma} \leftarrow (C_{m_i}, C_{\sigma,i}, pk_A)$ 
9:    $\pi_{\sigma} \leftarrow \text{ZK.PROVE}(pp_{zk}, \text{public}_{\sigma}, \text{witness}_{\sigma})$ 
10:   $S_{B,0} \leftarrow \emptyset$ 
11:  for  $j = 1$  to  $N$  do
12:     $R_{B,i,j} \leftarrow (\text{seed}_i = \text{PRF}(\omega_{\sigma,i} \parallel \omega_{m,i} \parallel \kappa_2) \wedge$ 
13:       $v_{i,j} = \text{HASH}(\text{seed}_i \parallel j \parallel \kappa_1) \wedge$ 
14:       $B_{i,j} = \langle \Delta W_{i,j}, v_{i,j} \rangle$ 
15:     $\omega_{B,j} \leftarrow (\Delta W_{i,j}, B_{i,j}, \omega_{\Delta W,i}, \omega_{B,i}, \omega_{m,i}, \omega_{\sigma,i})$ 
16:     $S_{B,j} \leftarrow \text{ZK.FOLD}(S_{B,j-1}, R_{B,i,j}, \omega_{B,j})$ 
17:  end for
18:   $\pi_B^{rec} \leftarrow \text{ZK.PROVE}(S_{B,N})$ 
19:   $S_{\Delta,0} \leftarrow \emptyset$ 
20:  for  $j = 1$  to  $N$  do
21:     $R_{\Delta,i,j} \leftarrow (\Delta W_{i,j} = W_{i,j} - W_{0,j})$ 
22:     $\text{witness}_{\Delta,j} \leftarrow (W_{0,j}, W_{i,j}, \Delta W_{i,j}, \omega_{W_0}, \omega_{W,i}, \omega_{\Delta W,i})$ 
23:     $S_{\Delta,j} \leftarrow \text{ZK.FOLD}(S_{\Delta,j-1}, R_{\Delta,i,j}, \text{witness}_{\Delta,j})$ 
24:  end for
25:   $\pi_{\Delta}^{rec} \leftarrow \text{ZK.PROVE}(S_{\Delta,N})$ 
26:   $\text{witness}_{\tau} \leftarrow (m_i, \omega_{m_i})$ 
27:   $\text{public}_{\tau} \leftarrow (C_{m_i}, p, r, \tau)$ 
28:   $\pi_{\tau} \leftarrow \text{ZK.PROVE}(pp_{zk}, \text{public}_{\tau}, \text{witness}_{\tau})$ 
29:   $\pi \leftarrow (\pi_{\sigma}, \pi_B^{rec}, \pi_{\Delta}^{rec}, \pi_{\tau})$ 
30:  Return  $(r, \pi)$ 
31: end procedure
    
```

The recursive weight consistency verification ensures that the committed weight differences correctly represent the transformation from the base model to the model across all layers. The verifier executes $\{0, 1\} \leftarrow \text{ZK.VERIFY}(pp_{zk}, S_{\Delta,N}, \pi_{\Delta}^{rec})$

to validate that the relation $\Delta W_{i,j} = W_{i,j} - W_{0,j}$ holds for all layers $j \in \{1, \dots, N\}$. This verification prevents malicious providers from claiming that arbitrary weight differences correspond to legitimate fine-tuning processes.

The query transcript consistency verification links the entire provenance verification to the specific query-response interaction, ensuring that the proof corresponds to the actual user prompt and generated response rather than being a generic proof that could be reused across different interactions. The verifier checks

$$\{0, 1\} \leftarrow \text{ZK.VERIFY}(pp_{zk}, (C_{m_i}, p, r, \tau), \pi_{\tau})$$

To confirm that the query transcript τ was correctly computed using the committed metadata for the selected dataset, the user's prompt p , and the generated response r . This verification step binds the entire proof to the specific query context, preventing replay attacks where old proofs might be reused for different queries.

The final verification decision aggregates all individual verification results to determine whether the complete provenance chain

is valid. If all the verification steps are passed successfully, the \mathcal{U} or any verifier of their choice accepts the response and its claimed provenance. If any verification step fails, the entire proof is rejected, ensuring that partial or incomplete proofs cannot be accepted as valid. The verification process is summarized in Alg. 3.

Algorithm 3 Verification Phase of ZKPROV

```

1: procedure VERIFY( $p, r, \pi, C$ )
2:   Input: Prompt  $p$ , Response  $r$ , Proof  $\pi$ , Commitments  $C$ 
3:   Parse  $\pi \leftarrow (\pi_\sigma, \pi_B^{rec}, \pi_\Delta^{rec}, \pi_\tau)$ 
4:    $\tau \leftarrow \text{HASH}(C_{m,i} \parallel p \parallel r \parallel \kappa_3)$ 
5:    $\text{public}_\sigma \leftarrow (C_{m,i}, C_{\sigma,i}, pk_{C,\mathcal{A}})$ 
6:    $b_1 \leftarrow \text{ZK.VERIFY}(pp_{zk}, \text{public}_\sigma, \pi_\sigma)$ 
7:   Reconstruct  $S_{B,N}$  from public bindings
8:    $b_2 \leftarrow \text{ZK.VERIFY}(pp_{zk}, S_{B,N}, \pi_B^{rec})$ 
9:   Reconstruct  $S_{\Delta,N}$  from public weight consistencies
10:   $b_3 \leftarrow \text{ZK.VERIFY}(pp_{zk}, S_{\Delta,N}, \pi_\Delta^{rec})$ 
11:   $\text{public}_\tau \leftarrow (C_{m,i}, p, r, \tau)$ 
12:   $b_4 \leftarrow \text{ZK.VERIFY}(pp_{zk}, \text{public}_\tau, \pi_\tau)$ 
13:  if  $b_1 = 1 \wedge b_2 = 1 \wedge b_3 = 1 \wedge b_4 = 1$  then
14:    Return ACCEPT
15:  else
16:    Return REJECT
17:  end if
18: end procedure
    
```

The proposed design of the ZKPROV ensures that all necessary cryptographic structures are established during the LLM response generation’s offline phase, enabling efficient proof generation during query time while maintaining the security properties required for dataset provenance verification. In the next section, we define the security of our proposed scheme and formally prove that ZKPROV binds the model’s response to its authorized authenticated dataset.

5 Security and Privacy Analyses

This chapter provides a formal cryptographic analysis of the proposed ZKPROV framework introduced in Section 4. We define the security properties our scheme aims to achieve in Sec. 4.2, and provide its comprehensive proofs of security under standard cryptographic assumptions.

We analyze the security of our proposed scheme for a probabilistic polynomial-time (PPT) adversary \mathcal{A} whose resources, such as time and queries, are bounded by a polynomial in the security parameter λ . The provider can be a malicious adversary \mathcal{A}_P who has access to a signing oracle $\mathcal{O}_{\text{sign}}(\cdot)$, which returns σ_i . It can query the hash oracle $\mathcal{O}_{\text{hash}}(\cdot)$ polynomially many times, may deviate arbitrarily from the protocol specification, and attempts to generate valid proofs for responses derived from unauthorized datasets.

The users are considered semi-honest adversaries \mathcal{A}_U who follow the protocol steps correctly but attempt to learn unauthorized information about dataset contents or model parameters from proofs and public information. Additionally, it has access to all public protocol outputs.

THEOREM 1 (SECURITY OF ZKPROV). *The ZKPROV provides soundness, dataset exposure binding, and Zero-Knowledge (Privacy) if the*

underlying KZG commitment scheme is computationally binding and perfectly hiding, the BLS signature scheme is unforgeable, the hash function HASH is collision-resistant, and the HyperNova proof system provides computational soundness and computational zero-knowledge

The proof starts by stating that if there exists a PPT adversary \mathcal{A} who can break the soundness and privacy of ZKPROV, we can use this adversary to construct \mathcal{B} to break the corresponding underlying assumption.

We formally define the security experiments and explain the advantage notation used throughout our analysis, where advantages are negligible functions of λ . Therefore, we have $\text{Adv}_{\mathcal{B}_1}^{\text{BLS}}(\lambda)$: Advantage of adversary \mathcal{B}_1 in breaking BLS signature unforgeability, $\text{Adv}_{\mathcal{B}_2}^{\text{Bind}}(\lambda)$: Advantage of adversary \mathcal{B}_2 in breaking KZG commitment binding (binding values), $\text{Adv}_{\mathcal{B}_3}^{\text{Bind}}(\lambda)$: Advantage of adversary \mathcal{B}_3 in breaking KZG commitment binding (weight consistency), $\text{Adv}_{\mathcal{B}_6}^{\text{Hide}}(\lambda)$: Advantage of adversary \mathcal{B}_6 in breaking KZG commitment hiding, $\text{Adv}_{\mathcal{B}_5}^{\text{CR}}(\lambda)$: Advantage of adversary \mathcal{B}_5 in finding hash collisions, $\text{Adv}_{\mathcal{B}_4}^{\text{HN-Sound}}(\lambda)$: Advantage of adversary \mathcal{B}_4 in breaking HyperNova soundness, and $\text{Adv}_{\mathcal{B}_7}^{\text{HN-ZK}}(\lambda)$: Advantage of adversary \mathcal{B}_7 in breaking HyperNova zero-knowledge.

DEFINITION 1 (SOUNDNESS EXPERIMENT). *The soundness experiment $\text{Exp}_{\mathcal{A}_P}^{\text{ZKPROV-Sound}}(\lambda)$ is defined as follows:*

- (1) *Challenger generates $(\text{crs}, sk_{C,\mathcal{A}}, pk_{C,\mathcal{A}})$ where crs contains domain separators $\kappa_1, \kappa_2, \kappa_3$.*
- (2) *Challenger gives crs and $pk_{C,\mathcal{A}}$ to \mathcal{A}_P .*
- (3) *Initialize empty query lists $Q_{\text{sign}} = \emptyset$ and $Q_{\text{hash}} = \emptyset$.*
- (4) *\mathcal{A}_P gets access to oracles $\mathcal{O}_{\text{sign}}(\cdot)$ and $\mathcal{O}_{\text{hash}}(\cdot)$.*
- (5) *\mathcal{A}_P outputs (pub^*, π^*) .*
- (6) *Return 1 iff $\text{VERIFY}(pp, \text{pub}^*, \pi^*) = \text{accept}$ AND $\text{ForgedProvenance}(\mathcal{A}_P, Q_{\text{sign}}) = 1$.*

We define $\text{FORGEDPROVENANCE}(\mathcal{A}_P, Q_{\text{sign}}) = 1$ iff \mathcal{A}_P ’s proof contains a signature on metadata $m^* \notin Q_{\text{sign}}$.

THEOREM 2 (SOUNDNESS). *For any PPT adversary \mathcal{A}_P :*

$$\Pr[\text{Exp}_{\mathcal{A}_P}^{\text{ZKPROV-Sound}}(\lambda) = 1] \leq \text{Adv}_{\mathcal{B}_1}^{\text{BLS}}(\lambda) + \text{Adv}_{\mathcal{B}_2}^{\text{Bind}}(\lambda) + \text{Adv}_{\mathcal{B}_3}^{\text{Bind}}(\lambda) + \text{Adv}_{\mathcal{B}_5}^{\text{CR}}(\lambda) + \text{Adv}_{\mathcal{B}_4}^{\text{HN-Sound}}(\lambda) + \frac{N \cdot d}{|\mathbb{F}_p|} \quad (1)$$

PROOF SKETCH. We construct \mathcal{B}_1 that uses HyperNova’s knowledge extractor \mathcal{E} to obtain (m_k, σ_k) . We prove this by analyzing the following four forgery types.

Dataset Authentication Forgery: If \mathcal{A}_P produces a valid proof for metadata $m_k \notin Q_{\text{sign}}$, we construct \mathcal{B}_1 that uses HyperNova’s knowledge extractor \mathcal{E} to obtain (m_k, σ_k) from π_σ and outputs this as a BLS forgery.

Binding Value Forgery: If $B_{i,j} \neq \langle \Delta W_{i,j}, v_{i,j} \rangle$, either (a) commitments are not binding, so \mathcal{B}_2 breaks KZG binding, or (b) distinct polynomials $f_{\text{honest}}(x) = \langle \Delta W_{i,j}, x \rangle$ and $f_{\text{malicious}}(x)$ agree at random challenge $v_{i,j}$. By Schwartz-Zippel, case (b) occurs with probability $\leq d/|\mathbb{F}_p|$ per layer. Challenge unpredictability follows from $v_{i,j} = \text{HASH}(\kappa_1 \parallel C_{\omega_{\sigma,i}} \parallel C_{\omega_{m,i}} \parallel i \parallel j)$, using the commitment hiding of randomness $\omega_{\sigma,i}, \omega_{m,i}$.

Weight Consistency Forgery: If $\Delta W_{i,j} \neq W_{i,j} - W_{0,j}$, but the proof verifies, then either the commitment binding is broken (\mathcal{B}_3 breaks KZG binding) or a false linear relation has a valid proof (\mathcal{B}_4 breaks HyperNova soundness).

Query Transcript Forgery: If $\tau^* \neq \text{HASH}(C_{m_i} \parallel p^* \parallel r^* \parallel \kappa_3)$, but proof passes the verification, \mathcal{B}_5 can extract both claimed and correct transcript inputs to find hash collision.

Since breaking the bounds of these defined cases equals breaking the security assumption, we conclude the proof by contradiction. \square

DEFINITION 2 (PRIVACY EXPERIMENT). *The privacy experiment $\text{Exp}_{\mathcal{A}_U}^{\text{ZKPROV-Priv}}(\lambda)$ is defined as follows:*

- (1) Challenger generates $(\text{crs}, \text{sk}_{C\mathcal{A}}, \text{pk}_{C\mathcal{A}})$.
- (2) Challenger gives crs and $\text{pk}_{C\mathcal{A}}$ to \mathcal{A}_U .
- (3) \mathcal{A}_U chooses datasets D_0, D_1 with metadata m_0, m_1 , prompt p , hyperparameters H , base model W_0 .
- (4) Challenger signs $\sigma_0 \leftarrow \text{BLS.SIGN}(\text{sk}_{C\mathcal{A}}, m_0)$, $\sigma_1 \leftarrow \text{BLS.SIGN}(\text{sk}_{C\mathcal{A}}, m_1)$.
- (5) Challenger flips $b \xleftarrow{\$} \{0, 1\}$.
- (6) Challenger computes $W_b \leftarrow \text{FINETUNE}(D_b, W_0, H)$, $r_b \leftarrow \text{LLM}(W_b, p)$, runs SETUP to generate C, Ω , and $\pi_b \leftarrow \text{PROVE}(C, \Omega, p)$.
- (7) Challenger gives (r_b, π_b, C) to \mathcal{A}_U .
- (8) \mathcal{A}_U outputs guess b' .
- (9) Return 1 iff $b' = b$.

DEFINITION 3 (PRIVACY ADVANTAGE). *The advantage of adversary \mathcal{A}_U in the privacy experiment is:*

$$\text{Adv}_{\mathcal{A}_U}^{\text{ZKPROV-Priv}}(\lambda) = \left| \Pr[\text{Exp}_{\mathcal{A}_U}^{\text{ZKPROV-Priv}}(\lambda) = 1] - \frac{1}{2} \right|$$

THEOREM 3 (ZERO-KNOWLEDGE). *For any PPT adversary \mathcal{A}_U :*

$$\text{Adv}_{\mathcal{A}_U}^{\text{ZKPROV-Priv}}(\lambda) \leq \text{Adv}_{\mathcal{B}_6}^{\text{Hide}}(\lambda) + \text{Adv}_{\mathcal{B}_7}^{\text{HN-ZK}}(\lambda)$$

PROOF SKETCH. We prove zero-knowledge via a hybrid argument with three hybrid statements, where each transition is proven computationally indistinguishable based on reductions. We define the simulator \mathcal{S}_{HN} that can generate computationally indistinguishable proofs for the defined statements, and the knowledge extractor \mathcal{E} extracts witnesses from valid proofs with probability $1 - \text{negl}(\lambda)$.

Hybrid H_0 (Real Experiment): The challenger runs the real privacy experiment where:

- All commitments are real: $C_{m_b} = \text{C.COMMIT}(m_b, \omega_{m_b})$, $C_{\sigma_b} = \text{C.COMMIT}(\sigma_b, \omega_{\sigma_b})$, $C_{\Delta W_b} = \text{C.COMMIT}(\Delta W_b, \omega_{\Delta W_b})$, $C_{B_b} = \text{C.COMMIT}(B_b, \omega_{B_b})$, $C_{W_0} = \text{C.COMMIT}(W_0, \omega_{W_0})$
- The proof is real: $\pi_b = (\pi_\sigma, \pi_B^{\text{rec}}, \pi_\Delta^{\text{rec}}, \pi_\tau) \leftarrow \text{PROVE}(C, \Omega, p)$
- Output: (r_b, π_b, C) where $C = (C_{m_b}, C_{\sigma_b}, C_{\Delta W_b}, C_{B_b}, C_{W_0})$

Hybrid H_1 (Simulated Proofs with Real Commitments): The real proofs are replaced with simulated ones. We keep the real commitments:

- Keep all real commitments: $C = (C_{m_b}, C_{\sigma_b}, C_{\Delta W_b}, C_{B_b}, C_{W_0})$ as in H_0
- Generate simulated proof components:
 - $\pi_\sigma \leftarrow \mathcal{S}_{\text{HN}}(C_{m_b}, C_{\sigma_b}, \text{pk}_{C\mathcal{A}})$
 - $\pi_B^{\text{rec}} \leftarrow \mathcal{S}_{\text{HN}}$
 - $\pi_\Delta^{\text{rec}} \leftarrow \mathcal{S}_{\text{HN}}$

$$- \pi_\tau \leftarrow \mathcal{S}_{\text{HN}}(C_{m_b}, p, r_b, \tau)$$

- Output: (r_b, π_b, C) where $\pi_b = (\pi_\sigma, \pi_B^{\text{rec}}, \pi_\Delta^{\text{rec}}, \pi_\tau)$

Hybrid H_2 (Simulated Commitments and Proofs):

Replace real commitments with random group elements:

- Generate random commitments: $\tilde{C}_{m_b}, \tilde{C}_{\sigma_b}, \tilde{C}_{\Delta W_b}, \tilde{C}_{B_b}, \tilde{C}_{W_0} \xleftarrow{\$} \mathbb{G}_1$
- Generate simulated proofs consistent with random commitments: $\tilde{\pi}_b \leftarrow \mathcal{S}_{\text{HN}}(\tilde{C})$
- Output: $(r_b, \tilde{\pi}_b, \tilde{C})$ where $\tilde{C} = (\tilde{C}_{m_b}, \tilde{C}_{\sigma_b}, \tilde{C}_{\Delta W_b}, \tilde{C}_{B_b}, \tilde{C}_{W_0})$

To analyze the distinguishability argument, we continue as follows.

Claim. $H_0 \approx_c H_1$ with $|\Pr[\mathcal{A}_U(H_0) = 1] - \Pr[\mathcal{A}_U(H_1) = 1]| \leq \text{Adv}_{\mathcal{B}_7}^{\text{HN-ZK}}(\lambda)$

Proof. We construct reduction \mathcal{B}_7 that uses any distinguisher \mathcal{D} between H_0 and H_1 to break HyperNova zero-knowledge:

- \mathcal{B}_7 receives HyperNova public parameters and a statement-proof pair (stmt, π) from the HyperNova zero-knowledge challenger, where π is either real or simulated.
- \mathcal{B}_7 simulates the ZKPROV privacy experiment by generating real commitments using actual witnesses.
- \mathcal{B}_7 embeds the challenge proof π as one of the ZKPROV proof components (e.g., π_σ) and generates the remaining components accordingly (real if π is real, simulated if π is simulated).
- \mathcal{B}_7 gives the resulting experiment output to \mathcal{D} and forwards \mathcal{D} 's output to the HyperNova challenger.

The advantage of \mathcal{B}_7 in the HyperNova zero-knowledge experiment equals the distinguishing advantage between H_0 and H_1 .

Claim. $H_1 \approx_c H_2$ with $|\Pr[\mathcal{A}_U(H_1) = 1] - \Pr[\mathcal{A}_U(H_2) = 1]| \leq \text{Adv}_{\mathcal{B}_6}^{\text{Hide}}(\lambda)$

Proof. We construct adversary \mathcal{B}_6 that uses any distinguisher \mathcal{D} between H_1 and H_2 to break KZG commitment hiding:

- \mathcal{B}_6 receives KZG public parameters and a commitment C^* from the KZG hiding challenger, where C^* is either a commitment to a real value or a random group element.
- \mathcal{B}_6 generates the ZKPROV privacy experiment where one of the commitments (e.g., C_{m_b}) is set to the challenge commitment C^* .
- The remaining commitments are generated consistently: real commitments if C^* is real, random group elements if C^* is random.
- All proofs are generated using the HyperNova simulator \mathcal{S}_{HN} to maintain consistency across both cases.
- \mathcal{B}_6 gives the resulting experiment output to \mathcal{D} and forwards \mathcal{D} 's output to the KZG challenger.

The advantage of \mathcal{B}_6 in the KZG hiding experiment equals the distinguishing advantage between H_1 and H_2 .

Claim. $\Pr[\mathcal{A}_U(H_2) = 1] = 1/2$

Proof. In hybrid H_2 , all commitments are random group elements and all proofs are simulated. The only component that depends on the choice bit b is the response r_b . However, by the semantic similarity requirement in the privacy experiment (the adversary chooses datasets D_0, D_1 that should produce similar responses), and the fact that all cryptographic components are now independent of the

actual datasets, the adversary has no advantage in distinguishing between $b = 0$ and $b = 1$.

A direct result of the three claims states that

$$\begin{aligned} \text{Adv}_{\mathcal{A}_U}^{\text{ZKPROV-Priv}}(\lambda) &= \left| \Pr[\mathcal{A}_U(H_0) = 1] - \frac{1}{2} \right| \\ &= |\Pr[\mathcal{A}_U(H_0) = 1] - \Pr[\mathcal{A}_U(H_2) = 1]| \\ &\leq |\Pr[\mathcal{A}_U(H_0) = 1] - \Pr[\mathcal{A}_U(H_1) = 1]| \\ &\quad + |\Pr[\mathcal{A}_U(H_1) = 1] - \Pr[\mathcal{A}_U(H_2) = 1]| \\ &\leq \text{Adv}_{\mathcal{B}_7}^{\text{HN-ZK}}(\lambda) + \text{Adv}_{\mathcal{B}_6}^{\text{Hide}}(\lambda). \end{aligned}$$

Privacy Implications: This proof establishes that:

- **Dataset Content Privacy:** The weight differences $\Delta W_{i,j}$ that encode dataset-specific information are perfectly hidden by KZG commitments and never revealed through the zero-knowledge proofs.
- **Model Parameter Privacy:** The fine-tuning transformations captured in $\Delta W_{i,j}$ remain confidential, protecting proprietary model improvements.
- **Training Process Privacy:** The binding values $B_{i,j}$ and their computation process remain hidden, preventing adversaries from learning about internal training dynamics.

□

PROOF OF THEOREM 1. The theorem follows directly from Theorems 2 and 3. Since all underlying cryptographic primitives satisfy their respective security properties by assumption, all advantage terms are negligible. The Schwartz-Zippel term $\frac{N \cdot d}{|\mathbb{F}_p|}$ is negligible for appropriately chosen parameters where $|\mathbb{F}_p|$ is exponential in λ while N, d are polynomial in λ . Therefore, ZKPROV achieves both soundness (including dataset exposure binding) and zero-knowledge privacy under the stated cryptographic assumptions. □

6 Evaluation and Comparison

This section details the empirical evaluation of the ZKPROV protocol proposed in Section 4 and its comparison to the state-of-the-art protocols.

We implemented the complete protocol using a fine-tuned causal language model (LLaMA-3.1-8B [12]) from the PubMedQA [14] dataset to demonstrate its efficacy in proving model provenance with zero-knowledge proofs. Our dataset includes approximately 1,000 biomedical questions, each paired with curated long-form answers and corresponding binary final decisions (yes/no). Each example is structured with a question p , a set of contextual passages, and a long-form answer, which we denote as r , the LLM’s response. Each data sample is preprocessed into a format that combines the question and context as input and concatenates the final decision with the long-form answer as the expected output. The dataset is split into 800 training samples, 100 validation samples, and 100 held-out test samples.

The authorized hyperparameter tuple $H = (\eta, B, E, O)$, where $\eta = 5 \times 10^{-5}$ is the learning rate, $B = 8$ is the effective batch size (accumulated from microbatches of 2), $E = 3$ is the number of training epochs, and O denotes the AdamW optimizer.

Our implementation leverages Microsoft’s Nova [19] recursive proof system built on a cycle of BN254 and Grumpkin elliptic curves

with scalar field \mathbb{F}_p where $p \approx 2^{254}$, providing 128-bit computational security suitable for production deployment. All experiments are conducted on a MacBook Air with an Apple M4 chip featuring 10 cores (4 performance and 6 efficiency cores) and 16 GB of unified memory.

We evaluated the computational effort of the \mathcal{P} to generate the complete proof $\pi = (\pi_\sigma, \pi_B^{\text{rec}}, \pi_\Delta^{\text{rec}}, \pi_\tau)$, described in Algorithm 2, and the \mathcal{U} ’s effort to verify each component as described in Algorithm 3 and summarize them in Table 1.

Table 1: ZKPROV’s Parties Benchmarks for 1 Dataset

| Layers | \mathcal{P} ’s Time (ms) | \mathcal{U} ’s Time (ms) |
|--------|----------------------------|----------------------------|
| 8 | 843.06 | 427.64 |
| 16 | 1,128.70 | 727.44 |
| 32 | 1,450.20 | 1,055.40 |

The results in Table 1 highlight ZKPROV’s sublinear scaling, with the prover’s cost growing as $O(n \log n)$, where n is the number of transformer layers. This is an improvement over the naive approaches, achieved via recursive folding mechanism that reduces cryptographic overhead.

For instance, scaling from 8 to 32 layers (4× increase) results in a prover time increase from 843.06 ms to 1,450.20 ms (1.72× growth) and verifier time growth from 427.64 ms to 1,055.40 ms (2.47× growth). The slightly superlinear verifier scaling reflects recursive validation accumulation. Sub-second proof generation and verification for up to 32 layers confirm ZKPROV’s real-time applicability, especially in critical domains like healthcare.

To assess multi-source performance, ZKPROV was benchmarked with up to 3 datasets, each comprising 800k parameters across 8 layers, representing realistic fine-tuning scenarios for specialized domains.

Table 2: ZKPROV Scalability with Datasets of 8 Layers

| # (D_i) | Constraints | \mathcal{P} ’s Time (ms) | \mathcal{U} ’s Time (ms) |
|-------------|-------------|----------------------------|----------------------------|
| 1 | ≈ 10k | 843.06 | 427.64 |
| 2 | ≈ 19k | 1,128.70 | 727.44 |
| 3 | ≈ 27k | 1,450.20 | 1,055.40 |

The multi-dataset scalability analysis highlights computational bottlenecks in ZKPROV’s architecture. As datasets increase, the constraint count grows linearly (10k → 19k → 27k for 1 → 2 → 3 datasets), but the prover’s cost exhibits superlinear growth due to Nova’s recursive folding. Each additional dataset requires cryptographic composition with prior states, making folding the dominant bottleneck. This involves costly elliptic curve operations like multi-scalar multiplications and pairings, which compound with recursive levels.

ZKPROV performs well for single-dataset scenarios (843.06 ms prover time) and remains viable for small-scale multi-dataset use cases (2-3 datasets, < 1.5 s prover time). However, scalability to larger datasets is limited by exponential computational costs, suggesting the need for optimizations like batch processing or hierarchical proofs for enterprise applications.

Although ZKPROV addresses verifiable dataset provenance, Table 3 compares related zero-knowledge systems targeting machine

learning integrity, focusing on cryptographic primitives, model types, and proof generation/verification times.

Table 3: Comparison of ZKPROV with the State-of-the-Art

| System | ZK Scheme | Proof Gen | Verification |
|---------------|------------------|------------|--------------|
| zkLLM [28] | zk-SNARK | 620 s | 2.35 s |
| zkCNN [21] | IP + Poly Commit | 88.3 s | 59.3 ms |
| ZEN [8] | zk-SNARK | 119.5 s | 18.6 ms |
| ZKPROV (Ours) | zk-SNARK | 122,187 ms | 1,498 ms |

ZEN was evaluated on LeNet using the CIFAR-10 dataset, focusing on quantization-friendly zero-knowledge proofs for neural network inference. zkCNN was tested on both LeNet and VGG16 models, using CIFAR-10 and synthetic inputs, and emphasizes scalability using interactive proofs (IP), where a prover and verifier exchange messages to establish correctness, and polynomial commitments, which allow the prover to commit to polynomials and later reveal evaluations without revealing the full polynomial. zkLLM targets transformer-based language models like LLaMA-7B and introduces optimizations for attention layers, achieving efficient zk-SNARKs for large-scale models.

ZKPROV introduces a new class of verifiable ML protocols focused on dataset attribution, proving that a model’s response originates from a model trained on an authorized dataset. This is crucial in regulated areas like healthcare, where data compliance is prioritized over computational reproducibility. Although not directly comparable, we include this performance snapshot to show that ZKPROV maintains reasonable overhead in a privacy-constrained environment.

7 Discussion and Future Work

This section discusses the limitations and future work of the designed ZKPROV.

Integration with Inference and PoT. ZKPROV complements existing verification frameworks by focusing specifically on dataset provenance rather than computational correctness. A direct application involves integration with inference verification protocols such as TeleSparse [23] for inference verification, and proof-of-training systems such as Kaizen [1] we can achieve comprehensive assurance verification of the complete LLM pipeline and provide end-to-end verification: dataset provenance through ZKPROV, correct training execution through proof-of-training, and accurate inference through systems like TeleSparse.

Multi-Authority and Updatable Datasets. The BLS signature scheme supports multi-authority scenarios by aggregating multiple signatures into a single compact one, which is useful in healthcare where approval from entities like institutional review boards and data governance committees is required. This method proves that all authorities authenticated the dataset without revealing specific identities or increasing the proof size.

Reckle Trees enable efficient updates for dynamic datasets through batch proof capabilities. When modifications occur, such as adding patient records or updating protocols, only affected portions are updated, allowing incremental changes without recomputing the entire dataset commitment or retraining models.

Topic-Based Watermarking. Integrating watermarking techniques into our dataset provenance framework presents a promising

area for future research. While traditional methods focus on general text attribution, our framework emphasizes the need for *topic-based watermarking* [24], which has shown effectiveness in general and specific contexts [25]. This approach could encode thematic signals, distinguishing between oncology and cardiology-related training data within model outputs. By adding this auxiliary provenance layer, we can cryptographically prove that a model was trained on authorized datasets and trace a response’s topical source through statistical watermark detection. This is particularly beneficial in hybrid deployments, enhancing accountability by differentiating outputs influenced by fine-tuning from those shaped by retrieval-time context.

Combining Fine-Tuning and Retrieval Modes. In our work, ZKPROV currently focuses on proving dataset integrity for models. However, many real-world applications rely on retrieval-augmented generation (RAG), where an LLM dynamically pulls context from a hosted dataset at inference time. In future work, we aim to explore how ZKPROV can be extended to hybrid settings where certain training datasets are fully embedded into the model (via fine-tuning), and others are integrated via verifiable retrieval. For example, if multiple authorized datasets are structured under a common format, our system could prove that an output is derived from a combination of embedded knowledge and approved external sources, all under cryptographic provenance constraints. This allows fine-grained response verification to show that a model is authorized and that specific facts came from verifiably approved documents.

Enhancing Privacy and Efficiency. ZKPROV ensures the confidentiality of training datasets during verification but lacks formal protections against inference-based leakage from repeated query-response pairs. Introducing differential privacy could enhance defenses against statistical disclosure attacks by providing bounds on information leakage alongside zero-knowledge proofs. Incremental learning with evolving datasets presents cryptographic challenges, as all commitments need re-generation with changes; however, using cryptographic accumulators or updatable commitment schemes could enable efficient updates without restarting the proof process. Finally, developing formal verification tools for cryptographic protocols will bolster confidence in production, as the framework’s implementation complexity poses risks despite being secure under standard assumptions. Leveraging methods like proof-carrying code or verifiable compilation could address these challenges.

8 Conclusion

We introduced ZKPROV, a novel cryptographic framework that addresses the critical challenge of dataset provenance verification for LLM through zero-knowledge proofs, achieving a paradigm shift from computationally expensive complete training verification to practical statistical binding approaches. We deployed Nova’s recursive proof system with customized gates to prove the correctness of the statistical binding and significantly improve the efficiency of the verification process for data provenance for larger models. The framework provides formal security guarantees under standard cryptographic assumptions while addressing real-world regulatory compliance requirements in sensitive domains.

References

- [1] Kasra Abbaszadeh, Christodoulos Pappas, Jonathan Katz, and Dimitrios Papadopoulos. 2024. Zero-knowledge proofs of training for deep neural networks. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 4316–4330.
- [2] Peter Martey Addo, Dominique Guegan, and Bertrand Hassani. 2018. Credit risk analysis using machine and deep learning models. *Risks* 6, 2 (2018), 38.
- [3] John Armour and Mari Sako. 2020. AI-enabled business models in legal services: from traditional law firms to next-generation law companies? *Journal of Professions and Organization* 7, 1 (2020), 27–46.
- [4] Madhu Aswathy. 2025. Machine Learning for Autonomous Systems: Navigating Safety, Ethics, and Regulation In. (2025).
- [5] K Arun Bhavsar, Jimmy Singla, Yasser D Al-Otaibi, Oh-Young Song, Yousaf Bin Zikria, and Ali Kashif Bashir. 2021. Medical diagnosis using machine learning: a statistical review. *Computers, Materials and Continua* 67, 1 (2021), 107–125.
- [6] Dan Boneh, Ben Lynn, and Hovav Shacham. 2001. Short signatures from the Weil pairing. In *International conference on the theory and application of cryptography and information security*. Springer, 514–532.
- [7] Bing-Jyue Chen, Suppakit Waiwitlikhit, Ion Stoica, and Daniel Kang. 2024. Zkml: An optimizing system for ml inference in zero-knowledge proofs. In *Proceedings of the Nineteenth European Conference on Computer Systems*. 560–574.
- [8] Boyuan Feng, Lianke Qin, Zhenfei Zhang, Yufei Ding, and Shumo Chu. 2021. ZEN: Efficient Zero-Knowledge Proofs for Neural Networks. *IACR Cryptol. ePrint Arch.* 2021 (2021), 87.
- [9] Boyuan Feng, Zheng Wang, Yuke Wang, Shu Yang, and Yufei Ding. 2024. ZENO: A Type-based Optimization Framework for Zero Knowledge Neural Network Inference. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*. 450–464. Reference [9] in zkPyTorch.pdf.
- [10] Sanjam Garg, Aarushi Goel, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoudy, Guru-Vamsi Policharla, and Mingyuan Wang. 2023. Experimenting with zero-knowledge proofs of training. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*. 1880–1894.
- [11] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. 2016. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *International conference on machine learning*. PMLR, 201–210.
- [12] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [13] Meng Hao, Hanxiao Chen, Hongwei Li, Chenkai Weng, Yuan Zhang, Haomiao Yang, and Tianwei Zhang. 2024. Scalable zero-knowledge proofs for non-linear functions in machine learning. In *33rd USENIX Security Symposium (USENIX Security 24)*. 3819–3836.
- [14] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146* (2019).
- [15] Daniel Kang, Tatsunori Hashimoto, Ion Stoica, and Yi Sun. 2022. Scaling up trustless DNN inference with zero-knowledge proofs. *arXiv preprint arXiv:2210.08674* (2022).
- [16] Daniel Kang, Tatsunori Hashimoto, Ion Stoica, and Yi Sun. 2023. Scaling up Trustless DNN Inference with Zero-Knowledge Proofs. In *NeurIPS 2023 Workshop on Regulatable Machine Learning*. <https://nips.cc/virtual/2023/80626> Reference [15] in zkPyTorch.pdf.
- [17] Aniket Kate, Gregory M Zaverucha, and Ian Goldberg. 2010. Constant-size commitments to polynomials and their applications. In *International conference on the theory and application of cryptography and information security*. Springer, 177–194.
- [18] Abhiram Kothapalli and Srinath Setty. 2024. HyperNova: Recursive arguments for customizable constraint systems. In *Annual International Cryptology Conference*. Springer, 345–379.
- [19] Abhiram Kothapalli, Srinath Setty, and Ioanna Tzialla. 2022. Nova: Recursive zero-knowledge arguments from folding schemes. In *Annual International Cryptology Conference*. Springer, 359–388.
- [20] Bo Li, Peng Qi, Bo Liu, Shuai Di, Jingen Liu, Jiquan Pei, Jinfeng Yi, and Bowen Zhou. 2023. Trustworthy AI: From principles to practices. *Comput. Surveys* 55, 9 (2023), 1–46.
- [21] Tianyi Liu, Xiang Xie, and Yupeng Zhang. 2021. zkCNN: Zero Knowledge Proofs for Convolutional Neural Network Predictions and Accuracy. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security (CCS '21)*. ACM, 2968–2985. <https://doi.org/10.1145/3460120.3485379> Reference [16] in zkPyTorch.pdf, also cited as [17].
- [22] Tao Lu, Haoyu Wang, Wenjie Qu, Zonghui Wang, Jinye He, Tianyang Tao, Wenzhi Chen, and Jiaheng Zhang. 2024. An efficient and extensible zero-knowledge proof framework for neural networks. *Cryptology ePrint Archive* (2024).
- [23] Mohammad M Maheri, Hamed Haddadi, and Alex Davidson. 2025. TeleSparse: Practical Privacy-Preserving Verification of Deep Neural Networks. *Proceedings on Privacy Enhancing Technologies (PETS)* (2025).
- [24] Alexander Nemecek, Yuzhou Jiang, and Erman Ayday. 2024. Topic-based watermarks for LLM-generated text. *arXiv preprint arXiv:2404.02138* (2024).
- [25] Alexander Nemecek, Yuzhou Jiang, and Erman Ayday. 2025. The Feasibility of Topic-Based Watermarking on Academic Peer Reviews. *arXiv preprint arXiv:2505.21636* (2025).
- [26] Charalampos Papamanthou, Shravan Srinivasan, Nicolas Gailly, Ismael Hishon-Rezaizadeh, Andrus Salumets, and Stjepan Golemac. 2024. Reckle trees: Updatable merkle batch proofs with applications. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 1538–1551.
- [27] Jacob T Schwartz. 1980. Fast probabilistic algorithms for verification of polynomial identities. *Journal of the ACM (JACM)* 27, 4 (1980), 701–717.
- [28] Haochen Sun, Jason Li, and Hongyang Zhang. 2024. zkllm: Zero knowledge proofs for large language models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 4405–4419.
- [29] Florian Tramer and Dan Boneh. 2018. Slalom: Fast, verifiable and private execution of neural networks in trusted hardware. *arXiv preprint arXiv:1806.03287* (2018).
- [30] Zhibo Xing, Zijian Zhang, Jiamou Liu, Ziang Zhang, Meng Li, Liehuang Zhu, and Giovanni Russello. 2023. Zero-knowledge proof meets machine learning in verifiability: A survey. *arXiv preprint arXiv:2310.14848* (2023).
- [31] Jiaheng Zhang, Zhiyong Fang, Yupeng Zhang, and Dawn Song. 2020. Zero Knowledge Proofs for Decision Tree Predictions and Accuracy. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security (CCS '20)*. Association for Computing Machinery, New York, NY, USA, 2039–2053. <https://doi.org/10.1145/3372297.3417278> Reference [26] in zkPyTorch.pdf.
- [32] Jiaheng Zhang, Zhiyong Fang, Yupeng Zhang, and Dawn Song. 2020. Zero knowledge proofs for decision tree predictions and accuracy. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*. 2039–2053.
- [33] Lingchen Zhao, Qian Wang, Cong Wang, Qi Li, Chao Shen, and Bo Feng. 2021. Veriml: Enabling integrity assurances and fair payments for machine learning as a service. *IEEE Transactions on Parallel and Distributed Systems* 32, 10 (2021), 2524–2540.