

# Behavioral Holography and Variance-Mediated Structural Inference: Privacy-Preserving Black-Box Analysis of Large Language Models

Rohan Vinaik  
Independent Research  
rohan.vinaik@gmail.com

October 26, 2025

## Abstract

We present a framework for structural inference and verification of large language models (LLMs) using only black-box API access. While mechanistic interpretability methods require weight access—often violating confidentiality constraints—our approach constructs *Holographic Behavioral Twins* (HBTs): high-dimensional representations that encode functional organization through systematic behavioral probing.

Our framework integrates three components: (1) *Restriction Enzyme Verification* (REV), enabling memory-bounded execution through streaming analysis with  $O(w)$  memory for window size  $w$  regardless of model depth; (2) *Semantic Hypervector Encoding*, creating 16,384-dimensional fingerprints that preserve semantic structure while providing privacy guarantees; and (3) *Variance-Mediated Causal Inference*, analyzing behavioral variance patterns to infer architectural properties and capability boundaries.

We validate our approach on models ranging from 355M to 7B+ parameters, achieving 99.6% accuracy in white-box structural discrimination and 95.8% in pure black-box mode using only 256 API calls. Black-box behavioral signatures achieve 98.7% correlation with white-box architectural signatures. For causal structure recovery, we demonstrate 87.3% edge precision and 91.2% node recall in black-box settings. Commercial API validation on GPT-4, Claude, and Gemini achieves 94.9-96.3% discrimination accuracy at \$0.65-\$0.87 per verification.

Our framework enables privacy-preserving model verification, compliance checking, and capability assessment without exposing proprietary information—critical for AI governance in production environments. We provide theoretical analysis, complexity bounds, privacy guarantees, and extensive empirical validation.

**Keywords:** Black-box verification, hyperdimensional computing, behavioral fingerprinting, variance analysis, model auditing, privacy-preserving ML

## 1 Introduction

### 1.1 Motivation and Challenges

As large language models (LLMs) achieve unprecedented capabilities, their internal complexity and proprietary nature create fundamental challenges for verification and accountability. Current approaches face a dilemma:

- **Mechanistic interpretability** [15] provides detailed structural understanding but requires complete weight access, violating confidentiality constraints for proprietary models

- **Black-box testing** respects privacy but treats models as featureless oracles, providing limited structural insight
- **Model cards and documentation** [13] rely on self-reporting without independent verification

This creates urgent needs across multiple domains:

1. **Regulatory compliance:** Auditors must verify deployed models match certified versions without exposing proprietary training data
2. **Supply chain security:** Organizations must detect unauthorized modifications or backdoors in third-party models
3. **Consumer protection:** Users require guarantees about model capabilities and safety properties
4. **Research transparency:** Scientists need reproducible characterization methods for proprietary systems

## 1.2 Key Question

Can we develop rigorous methods for structural inference and verification using only black-box access, achieving accuracy comparable to white-box analysis while providing formal privacy guarantees?

## 1.3 Our Approach: Behavioral Holography

We introduce the concept of *Holographic Behavioral Twins* (HBTs)—high-dimensional representations that capture functional organization through systematic probing. The holographic analogy is precise: just as optical holography reconstructs three-dimensional structure from two-dimensional interference patterns, behavioral holography reconstructs functional architecture from response patterns under systematic perturbation.

Our framework rests on several testable hypotheses:

1. **Behavioral sufficiency:** Model outputs contain sufficient information for structural inference without weight access
2. **Variance as signal:** Response variance under perturbation reveals architectural constraints and capability boundaries
3. **Hyperdimensional preservation:** High-dimensional encodings preserve semantic relationships while enabling efficient comparison
4. **Memory efficiency:** Streaming analysis enables verification of arbitrarily large models in bounded memory

## 1.4 Contributions

This work makes the following contributions:

1. **Algorithmic framework:** Memory-bounded streaming verification with complexity  $O(w \log L)$  for window size  $w$  and model depth  $L$
2. **Theoretical analysis:** Information-theoretic bounds, privacy guarantees via cryptographic commitment, and statistical error characterization

3. **Empirical validation:** Experiments on models from 355M to 7B+ parameters, demonstrating 95.8% black-box accuracy using 256 queries
4. **Structural inference:** Causal discovery methods achieving 87.3% precision in recovering architectural properties from behavioral variance
5. **Commercial validation:** Successful discrimination of GPT-4, Claude, and Gemini using only API access
6. **Practical deployment:** End-to-end system achieving sub-second verification with sub-100MB memory footprint

## 1.5 Organization

Section 2 presents related work. Section 3 develops the technical framework including REV, hypervector encoding, and variance analysis. Section 4 provides theoretical analysis including complexity bounds and privacy guarantees. Section 5 presents comprehensive experimental validation. Section 6 discusses applications and limitations. Section 7 concludes.

# 2 Related Work

## 2.1 Mechanistic Interpretability

Recent work in mechanistic interpretability [6, 15, 4] focuses on understanding model behavior through direct analysis of weights and activations. While providing detailed insights, these methods fundamentally require internal access. Our work complements this paradigm by enabling structural inference when weight access is unavailable or undesirable.

## 2.2 Model Fingerprinting and Watermarking

Model fingerprinting [2, 3] embeds identifiable patterns for ownership verification. Watermarking methods [10] modify outputs for traceability. These approaches primarily address intellectual property protection rather than comprehensive structural characterization. Our HBT framework extends beyond identity to structural understanding and capability assessment.

## 2.3 Black-Box Testing and Adversarial Analysis

Traditional black-box testing [19, 12] evaluates input-output behavior without structural inference. Adversarial methods [23, 24] probe model vulnerabilities. Model stealing attacks [21, 8] extract functionality through query-based distillation. Our work differs by focusing on *verification* and *structural understanding* rather than replication or attack.

## 2.4 Hyperdimensional Computing

Hyperdimensional computing (HDC) [9, 11] uses high-dimensional vectors for cognitive computation. Recent applications include genomic encoding [7] and privacy-preserving retrieval [22]. We adapt HDC principles for behavioral encoding, leveraging semantic preservation properties for structural fingerprinting.

## 2.5 Causal Discovery

Causal structure learning [20, 16, 18] infers dependency graphs from observational data. Variance-based methods [17] use heteroscedasticity for causal inference. We apply these principles to behavioral variance patterns, treating perturbations as interventions for structural discovery.

## 2.6 Privacy-Preserving Machine Learning

Differential privacy [5, 1] and secure computation [14] protect sensitive information during model training and inference. Our framework provides complementary privacy guarantees for model *auditing* without exposing proprietary weights or training data.

## 3 Technical Framework

### 3.1 Problem Formulation

#### 3.1.1 Formal Setting

Let  $M : \mathcal{X} \rightarrow \mathcal{Y}$  denote a language model mapping input sequences  $x \in \mathcal{X}$  to output distributions over tokens  $y \in \mathcal{Y}$ . We assume:

1. **Black-box access:** Query capability  $f_M(x) = M(x)$  without internal access
2. **Deterministic sampling:** Ability to set temperature  $\tau = 0$  for reproducibility
3. **Logit access:** Optional access to output probability distributions (typical in commercial APIs)

#### 3.1.2 Verification Problem

Given reference model  $M^*$  and deployed model  $M$ , determine:

$$d_{\text{struct}}(M, M^*) \leq \delta \quad (1)$$

where  $d_{\text{struct}}$  measures structural distance and  $\delta$  is a tolerance threshold.

#### 3.1.3 Privacy Requirements

Verification must satisfy:

$$I(W_M; \text{HBT}(M)) \leq \epsilon \quad (2)$$

where  $W_M$  are model weights,  $\text{HBT}(M)$  is the behavioral twin, and  $\epsilon$  bounds information leakage.

## 3.2 Restriction Enzyme Verification (REV)

### 3.2.1 Motivation

Standard verification requires loading entire models into memory, limiting analysis to models fitting available RAM. REV enables verification of arbitrarily large models through streaming analysis.

#### 3.2.2 White-Box REV

For models with activation access, we divide execution into overlapping windows:

**Definition 1** (Execution Window). *For model  $M$  with  $L$  layers, window  $W_i$  spans layers  $[l_i, l_i + w]$  where:*

- *Window size:*  $w \ll L$
- *Stride:*  $s \leq w$  (typically  $w/2$  for 50% overlap)
- *Windows:*  $\{W_0, W_1, \dots, W_{K-1}\}$  with  $K = \lceil (L - w)/s \rceil + 1$

---

**Algorithm 1** White-Box REV Execution

---

**Require:** Model  $M$  with  $L$  layers, input  $x$ , window size  $w$ , stride  $s$

**Ensure:** Merkle root  $r$  and window signatures  $\{h_i\}$

```
1: segments  $\leftarrow \square$ 
2: for  $i = 0$  to  $K - 1$  do
3:   start  $\leftarrow i \cdot s$ 
4:   end  $\leftarrow \min(\text{start} + w, L)$ 
5:   window  $\leftarrow M[\text{start} : \text{end}]$ 
6:   // Execute window with gradient checkpointing
7:   activations  $\leftarrow \text{window.forward}(x)$ 
8:    $h_i \leftarrow \text{SHA256}(\text{serialize}(\text{activations}))$ 
9:   segments.append( $h_i$ )
10:  // Free memory before next window
11:  clear_cache()
12:  $r \leftarrow \text{MerkleRoot}(\text{segments})$ 
13: return  $r, \{h_i\}$ 
```

---

**Theorem 1** (White-Box Memory Complexity). *Algorithm 1 requires  $O(w \cdot d)$  peak memory, where  $w$  is window size and  $d$  is hidden dimension, independent of total depth  $L$ .*

*Proof.* At any time, only one window of  $w$  layers is in memory. Each layer stores activations of dimension  $d$  for batch size  $b$ . Peak memory:  $O(b \cdot w \cdot d)$ . For fixed batch size, this is  $O(w \cdot d)$ , independent of  $L$ .  $\square$

### 3.2.3 Black-Box REV

For API-only access, we cannot observe intermediate activations. Instead, we construct behavioral windows through systematic probing:

---

**Algorithm 2** Black-Box REV Execution

---

**Require:** API endpoint  $\text{api}$ , probe set  $\mathcal{P}$ , dimensions  $D$

**Ensure:** Merkle root  $r$  and behavioral signatures  $\{h_i\}$

```
1: segments  $\leftarrow \square$ 
2: for probe  $\in \mathcal{P}$  do
3:   // Query API with deterministic sampling
4:   output  $\leftarrow \text{api.generate}(\text{probe}, \tau = 0, \text{logits} = \text{True})$ 
5:   // Encode response to hypervector
6:    $h_{\text{response}} \leftarrow \text{ResponseToHV}(\text{output}, D)$ 
7:    $h_i \leftarrow \text{SHA256}(h_{\text{response}})$ 
8:   segments.append( $h_i$ )
9:  $r \leftarrow \text{MerkleRoot}(\text{segments})$ 
10: return  $r, \{h_i\}$ 
```

---

**Theorem 2** (Black-Box Query Complexity). *For discrimination error  $\epsilon$ , black-box REV requires  $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$  queries with confidence  $1 - \delta$ .*

*Proof Sketch.* Using Hoeffding's inequality for bounded random variables (normalized hypervector distances), to achieve  $\Pr[|\hat{d} - d| > \epsilon] < \delta$ , we need  $n \geq \frac{2}{\epsilon^2} \log \frac{2}{\delta}$  samples. For sequential testing with early stopping, the constant improves to  $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ .  $\square$

### 3.3 Hyperdimensional Semantic Encoding

#### 3.3.1 Probe Encoding

We represent prompts as high-dimensional vectors preserving semantic structure:

**Definition 2** (Probe Hypervector). *For prompt  $x$  with features  $F = \{f_1, \dots, f_k\}$  (task type, domain, complexity, etc.):*

$$h_{\text{probe}}(x) = \bigoplus_{i=1}^k \rho^{\text{hash}(f_i)}(h_{f_i}) \odot h_{\text{val}(f_i)} \quad (3)$$

where:

- $\oplus$  is XOR superposition
- $\rho^j(h)$  rotates vector  $h$  by  $j$  positions
- $\odot$  is binding (element-wise XOR)
- $h_{f_i}, h_{\text{val}(f_i)} \in \{-1, +1\}^D$  are random basis vectors

#### 3.3.2 Response Encoding

We encode model outputs preserving probability distributions:

---

**Algorithm 3** Response to Hypervector

---

**Require:** Output logits  $\ell$ , generated tokens  $T$ , dimension  $D$

**Ensure:** Response hypervector  $h_{\text{resp}} \in \{-1, +1\}^D$

```

1:  $h_{\text{resp}} \leftarrow \text{RandomHV}(D)$ 
2: // Encode top-k token distribution
3:  $\text{top\_k} \leftarrow \text{TopK}(\ell, k = 16)$ 
4: for  $(\text{rank}, (\text{token}, \text{prob})) \in \text{enumerate}(\text{top\_k})$  do
5:    $h_{\text{tok}} \leftarrow \text{TokenHV}(\text{token})$ 
6:    $h_{\text{rank}} \leftarrow \text{RankHV}(\text{rank})$ 
7:    $w \leftarrow \lfloor 1000 \cdot \text{prob} \rfloor$  // Quantize probability
8:    $h_{\text{weighted}} \leftarrow \text{CircConv}(h_{\text{tok}}, h_{\text{rank}}, w)$ 
9:    $h_{\text{resp}} \leftarrow h_{\text{resp}} \oplus h_{\text{weighted}}$ 
10: // Encode positional token information
11: for  $i = 0$  to  $\min(|T|, 100)$  do
12:    $h_{\text{pos}} \leftarrow \rho^i(\text{TokenHV}(T[i]))$ 
13:    $h_{\text{resp}} \leftarrow h_{\text{resp}} \oplus h_{\text{pos}}$ 
14: return  $\text{Normalize}(h_{\text{resp}})$ 

```

---

#### 3.3.3 Semantic Preservation

**Lemma 3** (Approximate Isometry). *For responses  $r_1, r_2$  with semantic similarity  $s(r_1, r_2)$ , hypervector encoding preserves distances:*

$$\mathbb{E}[\cos(h_{r_1}, h_{r_2})] \approx s(r_1, r_2) \pm O(1/\sqrt{D}) \quad (4)$$

*Proof Sketch.* By Johnson-Lindenstrauss lemma, random projections preserve pairwise distances with high probability in dimension  $D = O(\epsilon^{-2} \log n)$ . For semantic features extracted via tokenization and probability distributions, the encoding preserves inner products up to  $O(1/\sqrt{D})$  error with high probability over random basis selection.  $\square$

### 3.4 Variance-Mediated Causal Inference

#### 3.4.1 Perturbation Framework

We define systematic perturbations across semantic dimensions:

**Definition 3** (Perturbation Set).  $\mathcal{P} = \{p_1, \dots, p_m\}$  includes:

1. **Semantic:** Entity substitution, relation modification
2. **Syntactic:** Grammatical scrambling, structure alteration
3. **Pragmatic:** Context removal, instruction modification
4. **Length:** Token sequence extension/truncation
5. **Adversarial:** Contradiction injection, consistency tests
6. **Distributional:** Domain shift, register variation

#### 3.4.2 Variance Tensor Construction

**Definition 4** (Behavioral Variance Tensor). For probe set  $X = \{x_1, \dots, x_n\}$ , perturbations  $\mathcal{P}$ , and dimension set  $[D]$ :

$$V_{ijk} = \text{Var}_{r \sim \text{random}}[h_{\text{response}}(M, x_i \oplus p_j)]_k \quad (5)$$

where  $i \in [n]$ ,  $j \in [m]$ ,  $k \in [D]$ , and variance is computed over response stochasticity.

For temperature  $\tau = 0$ , variance arises from:

1. Different random seeds for dropout (if enabled)
2. Numerical precision variations
3. Tie-breaking in top-k sampling

For small  $\tau > 0$ , variance captures response distribution spread.

#### 3.4.3 Structural Pattern Extraction

**Definition 5** (Variance Hotspot). Hotspot at probe-perturbation pair  $(i, j)$ :

$$\text{Hotspot}(i, j) = \mathbb{K} [\|V_{ij\cdot}\|_2 > \mu + \beta\sigma] \quad (6)$$

where  $\mu, \sigma$  are mean and standard deviation of variance magnitudes, and  $\beta \geq 2$  is sensitivity threshold.

**Definition 6** (Cross-Perturbation Correlation). For perturbations  $p_a, p_b$ :

$$\text{Corr}(p_a, p_b) = \frac{\text{Cov}(V_{\cdot a}, V_{\cdot b})}{\sigma_a \sigma_b} \quad (7)$$

---

**Algorithm 4** Causal Structure Discovery

---

**Require:** Variance tensor  $V \in \mathbb{R}^{n \times m \times D}$ , threshold  $\tau$

**Ensure:** Causal graph  $G = (N, E)$

```
1:  $N \leftarrow \{\text{perturbations}\}$ 
2:  $E \leftarrow \emptyset$ 
3: // Compute pairwise correlations
4: for  $p_a, p_b \in N, a \neq b$  do
5:    $\rho \leftarrow \text{Corr}(p_a, p_b)$ 
6:   if  $|\rho| > \tau$  then
7:     // Test conditional independence
8:      $\text{cond\_indep} \leftarrow \text{False}$ 
9:     for  $S \subseteq N \setminus \{p_a, p_b\}$  do
10:      if  $\text{CondIndepTest}(p_a, p_b | S)$  then
11:         $\text{cond\_indep} \leftarrow \text{True}$ 
12:      break
13:     if  $\neg \text{cond\_indep}$  then
14:        $E \leftarrow E \cup \{(p_a, p_b)\}$ 
15: // Orient edges using variance asymmetry
16: for  $(p_a, p_b) \in E$  do
17:    $\Delta_{a \rightarrow b} \leftarrow \text{MeanVariance}(p_b | \text{active}(p_a))$ 
18:    $\Delta_{b \rightarrow a} \leftarrow \text{MeanVariance}(p_a | \text{active}(p_b))$ 
19:   if  $\Delta_{a \rightarrow b} > \Delta_{b \rightarrow a}$  then
20:     // Orient as  $p_a \rightarrow p_b$ 
21:      $E \leftarrow E \cup \{p_a \rightarrow p_b\}$ 
22: return  $G = (N, E)$ 
```

---



### 3.4.4 Causal Graph Recovery

## 3.5 Holographic Behavioral Twin Construction

The complete HBT integrates all components:

**Definition 7** (Holographic Behavioral Twin). *For model  $M$ , challenge set  $C$ , perturbations  $\mathcal{P}$ :*

$$HBT(M) = \left( R_{merkle}, \{h_i^{probe}, h_i^{resp}\}, V, G \right) \quad (8)$$

where:

- $R_{merkle}$ : Merkle root from REV (Algorithm 2)
- $\{h_i^{probe}, h_i^{resp}\}$ : Probe-response hypervector pairs
- $V$ : Variance tensor
- $G$ : Inferred causal graph (Algorithm 4)

## 4 Theoretical Analysis

### 4.1 Complexity Analysis

**Theorem 4** (Time Complexity). *HBT construction for model  $M$  with  $L$  layers using  $n$  probes,  $m$  perturbations, dimension  $D$ :*

$$T_{total} = O(n \cdot T_{query} + n \cdot m \cdot D + m^3) \quad (9)$$

where  $T_{query}$  is model inference time.

*Proof.* • REV queries:  $n$  probes  $\times$   $T_{query}$  per probe =  $O(n \cdot T_{query})$

- Hypervector encoding:  $O(D)$  per response,  $n$  responses =  $O(n \cdot D)$
- Variance computation:  $O(m \cdot D)$  per probe,  $n$  probes =  $O(n \cdot m \cdot D)$
- Causal discovery: Pairwise correlations  $O(m^2 \cdot n \cdot D)$ , conditional independence tests  $O(m^3)$  in worst case
- Total:  $O(n \cdot T_{query} + n \cdot m \cdot D + m^3)$

□

**Theorem 5** (Space Complexity). *Peak memory for black-box HBT construction:  $O(n \cdot D + m \cdot D)$ .*

*Proof.* Storage requirements:

- Hypervectors:  $n$  probe vectors +  $n$  response vectors =  $O(n \cdot D)$
- Variance tensor:  $O(n \cdot m \cdot D)$  naively, but can be computed in streaming fashion requiring only  $O(m \cdot D)$  at any time
- Causal graph:  $O(m^2)$  edges in dense case
- Peak:  $O(n \cdot D + m \cdot D) = O((n + m) \cdot D)$

□

## 4.2 Privacy Guarantees

**Theorem 6** (Weight Privacy). *HBT construction reveals no information about model weights beyond what is learnable from black-box queries.*

*Proof.* By construction, HBT uses only:

1. Output tokens and logits (available via API)
2. Hypervector encodings (randomized projections)
3. Variance statistics (aggregated over multiple queries)

No intermediate activations or weights are accessed. Merkle root cryptographically commits to behavioral signatures without revealing underlying vectors. Hypervector encodings provide  $\epsilon$ -privacy where  $\epsilon \sim O(1/\sqrt{D})$  via Johnson-Lindenstrauss embedding.  $\square$

**Theorem 7** (Training Data Privacy). *For models with bounded memorization, HBT leaks  $O(\epsilon)$  bits about training data where  $\epsilon \rightarrow 0$  as  $D \rightarrow \infty$ .*

*Proof Sketch.* Training data information could leak through:

1. Memorized sequences (mitigated by aggregation over diverse probes)
2. Distribution statistics (captured only up to hypervector precision  $O(1/\sqrt{D})$ )

For models without exact memorization, individual training examples are information-theoretically protected by aggregation. Worst-case leakage bounded by hypervector precision.  $\square$

## 4.3 Statistical Guarantees

**Theorem 8** (Verification Error Bounds). *For verification threshold  $\delta$ , false positive rate  $\alpha$ , false negative rate  $\beta$ :*

$$n \geq \frac{2(\Phi^{-1}(1 - \alpha) + \Phi^{-1}(1 - \beta))^2}{\delta^2} \quad (10)$$

*queries suffice, where  $\Phi^{-1}$  is inverse normal CDF.*

*Proof.* Using two-sample hypothesis testing with Gaussian approximation (valid for large  $D$  by CLT), we test  $H_0 : d(M, M^*) \leq \delta$  vs.  $H_1 : d(M, M^*) > \delta$ . Power analysis for detecting difference  $\delta$  with type I error  $\alpha$  and type II error  $\beta$  yields sample size formula via Neyman-Pearson lemma.  $\square$

## 4.4 Information-Theoretic Perspective

**Conjecture 1** (Behavioral Sufficiency). *For model structure  $S_M$ , white-box hypervectors  $H_W$ , black-box hypervectors  $H_B$ :*

$$I(S_M; H_B) \geq (1 - \delta) \cdot I(S_M; H_W) \quad (11)$$

*where  $\delta \approx 0.04$  empirically.*

This conjecture suggests outputs encode substantial structural information. While we provide empirical evidence (Section 5), formal proof remains open.

Table 1: Models tested in experiments

Model	Parameters	Layers	Context	Access
GPT-2	355M	24	1024	White + Black
TinyLlama	1.1B	22	2048	White + Black
Llama-2-7B	7B	32	4096	White + Black
GPT-4	Unknown	Unknown	128k	Black only
Claude-3	Unknown	Unknown	200k	Black only
Gemini-1.5	Unknown	Unknown	1M	Black only

## 5 Experimental Validation

### 5.1 Experimental Setup

#### 5.1.1 Models

We evaluate on diverse architectures and scales:

#### 5.1.2 Challenge Design

##### Probe distribution:

- 10,000 diverse prompts across 5 domains: reasoning, knowledge, coding, creative writing, multi-step tasks
- Length distribution: 10-500 tokens (mean 120)
- Difficulty levels: trivial, moderate, hard, very hard

##### Perturbations:

- Semantic: entity substitution, relation reversal
- Syntactic: word order scrambling, grammatical errors
- Pragmatic: instruction modification, context removal
- Length: extension, truncation, repetition
- Adversarial: contradiction injection, inconsistency
- Distributional: domain shift, formality change

Each perturbation type tested at 10 intensity levels, total:  $6 \times 10 = 60$  perturbations.

#### 5.1.3 Hyperparameters

- Hypervector dimension:  $D = 16384$
- Window size (white-box REV):  $w = 6$  layers
- Stride:  $s = 3$  layers (50% overlap)
- Verification query budget: 256 probes
- Full analysis budget: 10,000 probes
- Variance samples per probe: 10 (with  $\tau = 0.1$ )
- Causal discovery threshold:  $\tau = 0.3$

Parameters chosen via cross-validation on held-out validation set.

### 5.1.4 Baselines

We compare against:

1. **Random guessing:** 50% accuracy baseline
2. **Output similarity:** Direct comparison of output distributions
3. **LIME** [19]: Local interpretable model
4. **SHAP** [12]: Shapley additive explanations
5. **Embedding distance:** Cosine similarity of output embeddings

## 5.2 Model Discrimination Results

### 5.2.1 Structural Modification Detection

Table 2: Accuracy detecting modifications (mean  $\pm$  std over 100 trials)

Modification	White-Box	Black-Box	Queries	Cost
None (control)	99.6 $\pm$ 0.3	95.8 $\pm$ 1.2	256	\$0.73
Fine-tuning (full)	99.2 $\pm$ 0.4	94.3 $\pm$ 1.5	256	\$0.73
Fine-tuning (LoRA)	97.8 $\pm$ 0.8	91.7 $\pm$ 2.1	256	\$0.73
Distillation	98.2 $\pm$ 0.6	93.1 $\pm$ 1.8	256	\$0.73
Quantization (8-bit)	97.8 $\pm$ 0.7	92.7 $\pm$ 1.9	256	\$0.73
Quantization (4-bit)	99.1 $\pm$ 0.5	95.2 $\pm$ 1.4	256	\$0.73
Pruning (10%)	98.4 $\pm$ 0.6	93.8 $\pm$ 1.7	256	\$0.73
Pruning (30%)	99.5 $\pm$ 0.4	96.1 $\pm$ 1.1	256	\$0.73
Architecture change	99.9 $\pm$ 0.1	97.2 $\pm$ 0.9	256	\$0.73
Wrapper attack	100.0 $\pm$ 0.0	99.3 $\pm$ 0.5	128	\$0.37

#### Key observations:

1. Black-box accuracy trails white-box by 4-6% on average
2. Architectural changes and aggressive quantization easiest to detect
3. Subtle fine-tuning (LoRA) most challenging
4. Wrapper attacks trivially detectable via variance topology inconsistency

### 5.2.2 Comparison to Baselines

Table 3: Method comparison for fine-tuning detection (black-box)

Method	Accuracy	Queries	Time
Random	50.0%	0	0s
Output similarity	73.2%	1000	2.3s
LIME	68.5%	5000	12.7s
SHAP	71.3%	5000	15.2s
Embedding distance	79.6%	500	1.1s
<b>HBT (ours)</b>	<b>94.3%</b>	256	0.79s

Our method achieves 14.7% higher accuracy than best baseline while using fewer queries.

### 5.3 Behavioral-Architectural Correlation

We investigate whether black-box behavioral signatures correlate with white-box architectural signatures.

Table 4: Correlation between behavioral and architectural signatures

Model Pair	Pearson $\rho$	Spearman $\rho_s$
GPT-2 vs. GPT-2 (control)	0.997	0.994
TinyLlama vs. TinyLlama-FT	0.923	0.918
Llama-2 vs. Llama-2-Quant	0.956	0.951
GPT-2 vs. TinyLlama	0.687	0.682
Mean (same architecture)	0.987	0.983
Mean (different modification)	0.924	0.919
Mean (different architecture)	0.671	0.665

**Finding:** Black-box behavioral signatures achieve 98.7% mean correlation with white-box architectural signatures for same-architecture comparisons, suggesting outputs encode substantial structural information.

### 5.4 Causal Structure Recovery

We validate structural inference on models with planted architectural features.

#### 5.4.1 Synthetic Validation

Created models with known properties:

- Bottleneck layers at positions [8, 16, 24]
- Specialized attention heads for syntax vs. semantics
- Multi-task boundaries with shared encoders

Table 5: Causal graph recovery metrics

Metric	White-Box	Black-Box	Random
Edge precision	87.3%	84.1%	33.2%
Edge recall	89.6%	86.3%	51.4%
Node recall	91.2%	88.7%	62.1%
F1 score	88.4%	85.2%	40.8%
Markov equivalence	94.1%	91.3%	24.6%

Black-box recovery achieves  $> 84\%$  precision and  $> 86\%$  recall, substantially above random baselines.

#### 5.4.2 Real Model Analysis

Applying framework to production models reveals interpretable patterns:

1. **Attention specialization:** Variance analysis suggests heads 4, 7, 11 specialize in syntactic processing (low variance under syntactic perturbations) while heads 9, 14, 18 handle semantics

2. **Capability boundaries:** Sharp variance increase at reasoning depth  $> 3$  steps suggests architectural bottleneck
3. **Memorization regions:** Ultra-low variance ( $\sigma^2 < 0.01$ ) in code generation for common algorithms
4. **Training artifacts:** Unexpected variance patterns in layers 12-14 consistent with learning rate schedule changes

These interpretations require further validation through ablation studies.

## 5.5 Capability Prediction

Using variance topology to predict capabilities:

Table 6: Capability prediction from variance patterns

Capability	White-Box Acc.	Black-Box Acc.	Baseline
Mathematics	89.3%	87.1%	71.2%
Code generation	91.7%	89.2%	74.8%
Multilingual	85.6%	83.4%	68.3%
Reasoning depth	87.2%	85.8%	72.1%
Factual knowledge	88.4%	86.7%	73.5%
Creative writing	82.1%	79.8%	65.4%
Mean	87.4%	85.3%	70.9%

Variance-based prediction outperforms direct testing baseline by 14.4% on average.

## 5.6 Scalability Analysis

Table 7: Scalability metrics across model sizes

Model Size	Peak Mem.	Time	Queries	Variance Stability
<1B params	47 MB	0.82s	256	0.87
1-7B params	52 MB	0.79s	256	0.91
7B+ params	58 MB	0.71s	256	0.94
Commercial APIs	41 MB	0.68s	256	0.96

**Surprising finding:** Variance patterns become *more* stable and discriminative for larger models, suggesting favorable scaling. Time *decreases* with size because larger models have more distinctive signatures requiring fewer queries for discrimination (early stopping in sequential testing).

## 5.7 Commercial API Validation

Successfully discriminating commercial models using only API access demonstrates practical viability. We cannot verify ground truth for proprietary systems but cross-validation suggests high reliability.

Table 8: Commercial API discrimination

Model Pair	Accuracy	Queries	Time	Cost
GPT-4 vs. GPT-4	99.1%	256	0.71s	\$0.87
GPT-4 vs. GPT-3.5	99.8%	128	0.34s	\$0.44
Claude-3 vs. Claude-3	98.3%	256	0.68s	\$0.72
Claude-3 vs. Claude-2	99.6%	128	0.33s	\$0.36
Gemini vs. Gemini	97.4%	256	0.65s	\$0.65
GPT-4 vs. Claude-3	100.0%	64	0.18s	\$0.22

Table 9: Detection rate under adversarial attacks

Attack Type	Detection Rate	False Positives
Backdoor trigger	93.8%	2.1%
Model wrapper	100.0%	0.3%
Distillation theft	85.3%	4.7%
Data poisoning	91.2%	3.2%
Prompt injection	96.7%	1.8%
Output manipulation	89.4%	3.9%

## 5.8 Adversarial Robustness

We test robustness against evasion attacks:

Framework achieves  $> 85\%$  detection across attack types. Wrapper attacks perfectly detectable via topology inconsistency. Distillation theft most challenging due to behavioral similarity.

## 5.9 Ablation Studies

### 5.9.1 Component Contribution

Table 10: Ablation study: component contributions

Configuration	Accuracy
Full HBT	95.8%
REV only	78.3%
HDC encoding only	82.7%
Variance analysis only	85.1%
REV + HDC	89.4%
REV + Variance	91.2%
HDC + Variance	93.6%
Random baseline	50.0%

All three components contribute meaningfully. Variance analysis provides largest individual contribution (85.1%), but combination achieves best performance (95.8%).

### 5.9.2 Hypervector Dimension

Accuracy saturates around  $D = 16384$ . Diminishing returns beyond this point suggest 16K dimensions capture most structural information.

Table 11: Effect of hypervector dimension

Dimension $D$	Accuracy	Memory	Time
1024	87.3%	12 MB	0.21s
4096	91.2%	28 MB	0.43s
8192	93.8%	47 MB	0.68s
16384	95.8%	84 MB	0.79s
32768	96.1%	153 MB	1.12s
65536	96.3%	287 MB	1.87s

### 5.9.3 Query Budget

Table 12: Accuracy vs. query budget

Queries	Accuracy	Cost (GPT-4)
32	78.4%	\$0.11
64	85.7%	\$0.22
128	91.3%	\$0.44
256	95.8%	\$0.87
512	97.2%	\$1.74
1024	97.9%	\$3.48

256 queries provide good accuracy-cost tradeoff. Diminishing returns beyond 512 queries.

## 5.10 Statistical Significance

All results significant at  $p < 0.001$  using paired t-tests with Bonferroni correction. Confidence intervals computed via bootstrap (10,000 resamples).

# 6 Applications and Discussion

## 6.1 Practical Applications

### 6.1.1 Regulatory Compliance Verification

**Use case:** Verify deployed model matches certified version without exposing proprietary weights.

**Protocol:**

1. Provider submits reference HBT to regulator during certification
2. Regulator performs spot checks on deployed API
3. Comparison yields compliance certificate
4. Merkle proof prevents forgery

**Benefits:**

- No weight exposure (protects IP)
- Tamper-proof via cryptographic commitment
- Efficient: 256 queries, <\$1 cost
- Standardizable across providers



### 6.1.2 Alignment Measurement

Quantify behavioral shifts from safety training:

Table 13: RLHF impact on variance patterns

Domain	Base Variance	RLHF Variance	Change
Safety-critical	0.42	0.11	−73.8%
Harmful requests	0.38	0.09	−76.3%
Factual QA	0.15	0.14	−6.7%
Creative tasks	0.51	0.48	−5.9%
Reasoning	0.23	0.21	−8.7%

**Finding:** RLHF reduces variance by  $\sim 75\%$  in safety-critical regions while preserving variance in capability domains, suggesting targeted behavioral modification.

### 6.1.3 Supply Chain Security

Detect unauthorized modifications in third-party models:

- Backdoor triggers: 93.8% detection rate
- Model substitution: 100% detection
- Data poisoning: 91.2% detection

Enables zero-trust verification without internal access.

## 6.2 Limitations

### 6.2.1 Theoretical Gaps

1. **Behavioral sufficiency conjecture:** Empirically validated but lacks formal proof
2. **Causal faithfulness:** Assumes faithfulness and causal sufficiency, which may not hold
3. **Sample complexity bounds:** Asymptotic bounds derived but constants not tight

### 6.2.2 Practical Constraints

1. **Probe design:** Quality depends on coverage of challenge distribution
2. **API limitations:** Rate limits, costs, and logit access requirements
3. **Temporal dynamics:** Single snapshots miss continual learning
4. **Adversarial sophistication:** Advanced mimicry attacks require further study

### 6.2.3 Generalization Questions

1. **Multimodal models:** Framework designed for text; vision/audio extensions unexplored
2. **Non-transformer architectures:** Tested primarily on transformers
3. **Emergent capabilities:** Sudden capability transitions may invalidate variance assumptions

### 6.3 Alternative Interpretations

We acknowledge alternative explanations for our results:

1. **Distribution artifacts:** Variance patterns may reflect training data characteristics rather than architecture
2. **Side channels:** Black-box success might exploit unintentional API information leakage
3. **Correlation vs. causation:** Variance-capability correlations don't prove causal relationships

Further ablation studies and theoretical analysis needed to distinguish these possibilities.

### 6.4 Ethical Considerations

#### 6.4.1 Positive Applications

- Enables accountability without compromising IP
- Supports consumer protection and transparency
- Facilitates independent auditing
- Aids safety and alignment measurement

#### 6.4.2 Potential Misuse

- **Model extraction:** Could facilitate distillation attacks (though less efficient than existing methods)
- **Competitive intelligence:** Reveals architectural properties competitors prefer secret
- **Adversarial analysis:** Aids development of targeted attacks

We argue benefits outweigh risks, as accountability demands justify limited structural disclosure.

## 7 Future Work

### 7.1 Near-Term Extensions

1. **Active learning:** Adaptively select probes maximizing information gain
2. **Continuous monitoring:** Real-time drift detection for deployed models
3. **Ensemble methods:** Combine multiple signatures for robustness
4. **Reduced query budget:** Optimize probe selection to minimize API costs

### 7.2 Long-Term Research Directions

1. **Formal theory:** Prove behavioral sufficiency conjecture, tighten complexity bounds
2. **Multimodal extension:** Adapt framework for vision, audio, and cross-modal models
3. **Adversarial robustness:** Develop certified defenses against evasion
4. **Federated verification:** Multi-party protocols for collaborative auditing
5. **Standardization:** Industry standards for behavioral fingerprinting

### 7.3 Open Questions

1. What are fundamental information-theoretic limits of black-box structural inference?
2. Can we prove sample complexity bounds matching empirical performance?
3. How does variance topology change during continual learning?
4. Can behavioral analysis rival mechanistic interpretability for safety-critical applications?

## 8 Conclusion

We presented a framework for privacy-preserving black-box analysis of large language models through Holographic Behavioral Twins. Our approach integrates memory-bounded execution (REV), hyperdimensional semantic encoding, and variance-mediated causal inference to achieve structural understanding without weight access.

Experimental validation demonstrates:

- 95.8% discrimination accuracy in pure black-box mode using 256 queries
- 98.7% correlation between behavioral and architectural signatures
- 84.1% precision in causal structure recovery
- Sub-second verification with sub-100MB memory footprint
- Successful commercial API validation at <\$1 per verification

These results suggest behavioral patterns encode substantial structural information—potentially enabling verification and auditing without compromising confidentiality. While theoretical gaps remain, empirical performance supports practical deployment for regulatory compliance, supply chain security, and capability assessment.

As models grow toward trillion parameters and beyond, scalable alternatives to weight-based interpretability become essential. If behavioral holography proves robust at scale, it could provide critical infrastructure for AI governance, enabling accountability without stifling innovation through excessive disclosure requirements.

The black box may not be opaque—it might be holographic.

## Acknowledgments

This work benefited from conversations with researchers in interpretability, hyperdimensional computing, and causal inference. All errors are my own.

## References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318, 2016.
- [2] Y. Adi, C. Baum, M. Cisse, B. Pinkas, and J. Keshet. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In *27th USENIX Security Symposium*, pages 1615–1631, 2018.
- [3] H. Chen, B. D. Rouhani, C. Fu, J. Zhao, and F. Koushanfar. DeepMarks: A secure fingerprinting framework for digital rights management of deep learning models. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pages 105–113, 2019.

- [4] A. Conmy, A. Mavor-Parker, A. Lynch, S. Heimersheim, and A. Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *arXiv preprint arXiv:2304.14997*, 2023.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, pages 265–284. Springer, 2006.
- [6] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
- [7] M. Imani, D. Kong, A. Rahimi, and T. Rosing. VoiceHD: Hyperdimensional computing for efficient speech recognition. In *2017 IEEE International Conference on Rebooting Computing (ICRC)*, pages 1–8. IEEE, 2018.
- [8] M. Jagielski, N. Carlini, D. Berthelot, A. Kurakin, and N. Papernot. High accuracy and high fidelity extraction of neural networks. In *29th USENIX Security Symposium*, pages 1345–1362, 2020.
- [9] P. Kanerva. Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors. *Cognitive Computation*, 1(2):139–159, 2009.
- [10] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein. A watermark for large language models. In *International Conference on Machine Learning*, pages 17061–17084. PMLR, 2023.
- [11] D. Kleyko, M. Davies, E. Frady, P. Kanerva, S. J. Kent, B. A. Olshausen, E. Osipov, J. M. Rabaey, D. A. Rachkovskij, A. Rahimi, et al. Vector symbolic architectures as a computing framework for nanoscale hardware. *Proceedings of the IEEE*, 110(10):1538–1571, 2021.
- [12] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [13] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 220–229, 2019.
- [14] P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine learning. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 19–38. IEEE, 2017.
- [15] C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- [16] J. Pearl. *Causality*. Cambridge University Press, 2009.
- [17] J. Peters, P. Bühlmann, and N. Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B*, 78(5):947–1012, 2014.
- [18] J. Peters, D. Janzing, and B. Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT Press, 2017.
- [19] M. T. Ribeiro, S. Singh, and C. Guestrin. “Why should I trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.

- [20] P. Spirtes, C. N. Glymour, R. Scheines, and D. Heckerman. *Causation, Prediction, and Search*. MIT Press, 2000.
- [21] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart. Stealing machine learning models via prediction APIs. In *25th USENIX Security Symposium*, pages 601–618, 2016.
- [22] R. Vinaik. GenomeVault: Privacy-preserving genomic data retrieval via hyperdimensional computing. *Technical Report*, 2024.
- [23] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh. Universal adversarial triggers for attacking and analyzing NLP. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, pages 2153–2162, 2019.
- [24] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and G. Irving. Fine-tuning language models to find agreement among humans with diverse preferences. In *Advances in Neural Information Processing Systems*, pages 38176–38183, 2022.