

# Executive Summary



# Where Lies Go to Hide: Probing the Geometry of LLM Hallucinations

Author - Arnav Raj ([arnavvraj.compsci@gmail.com](mailto:arnavvraj.compsci@gmail.com))  
Computer Science Dept, Indian Institute of Technology, Delhi

## The Problem

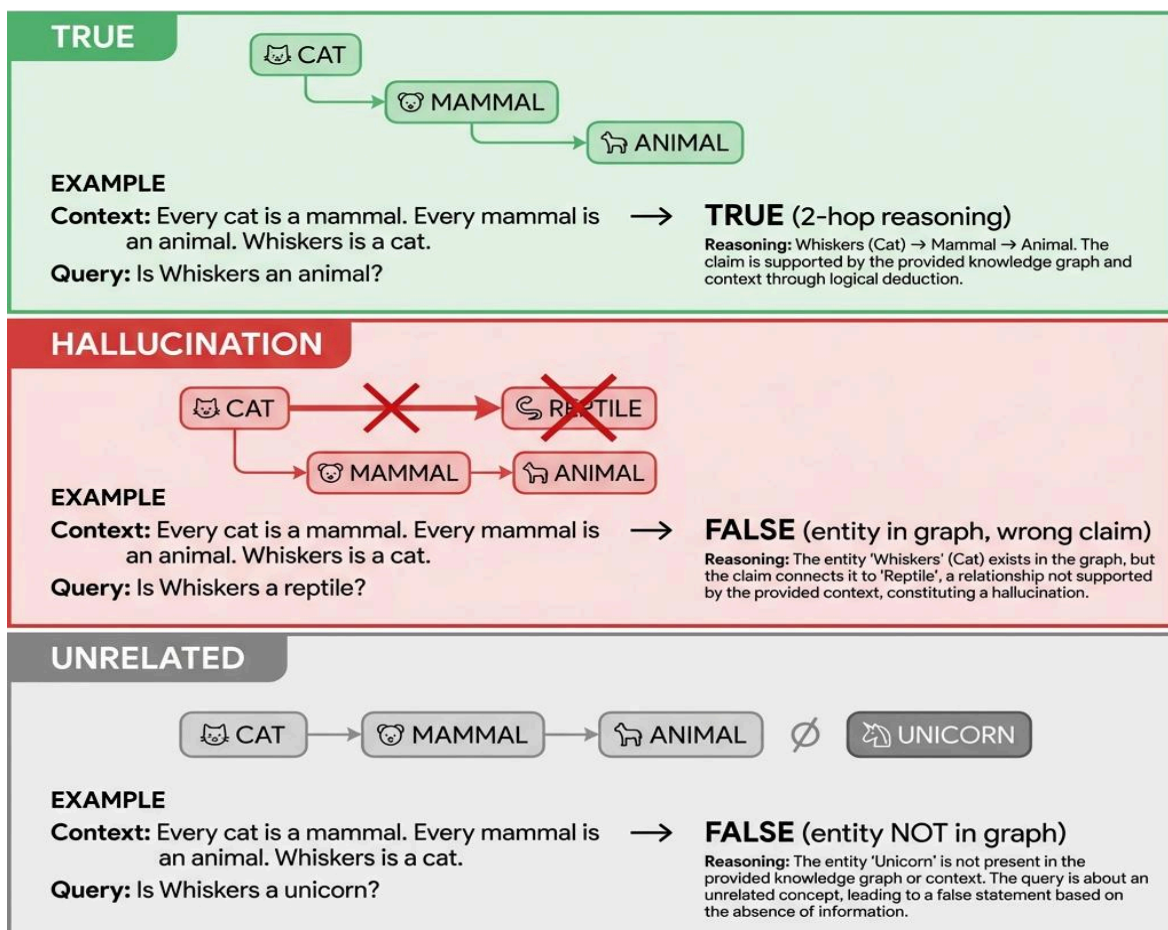
Large Language Models (LLMs) hallucinate—they produce confident, plausible-sounding statements that are factually false. Current detection methods require labeled hallucination datasets, which are expensive to create and may not generalize. We investigate whether the **geometry of LLM internal representations** can distinguish truth from hallucination *without requiring labeled hallucination training data*. [Repo Link](#)

## Core Question

**Can we detect hallucinations by analyzing where truthful vs fabricated outputs sit in the model's representation space?**

We hypothesize that:

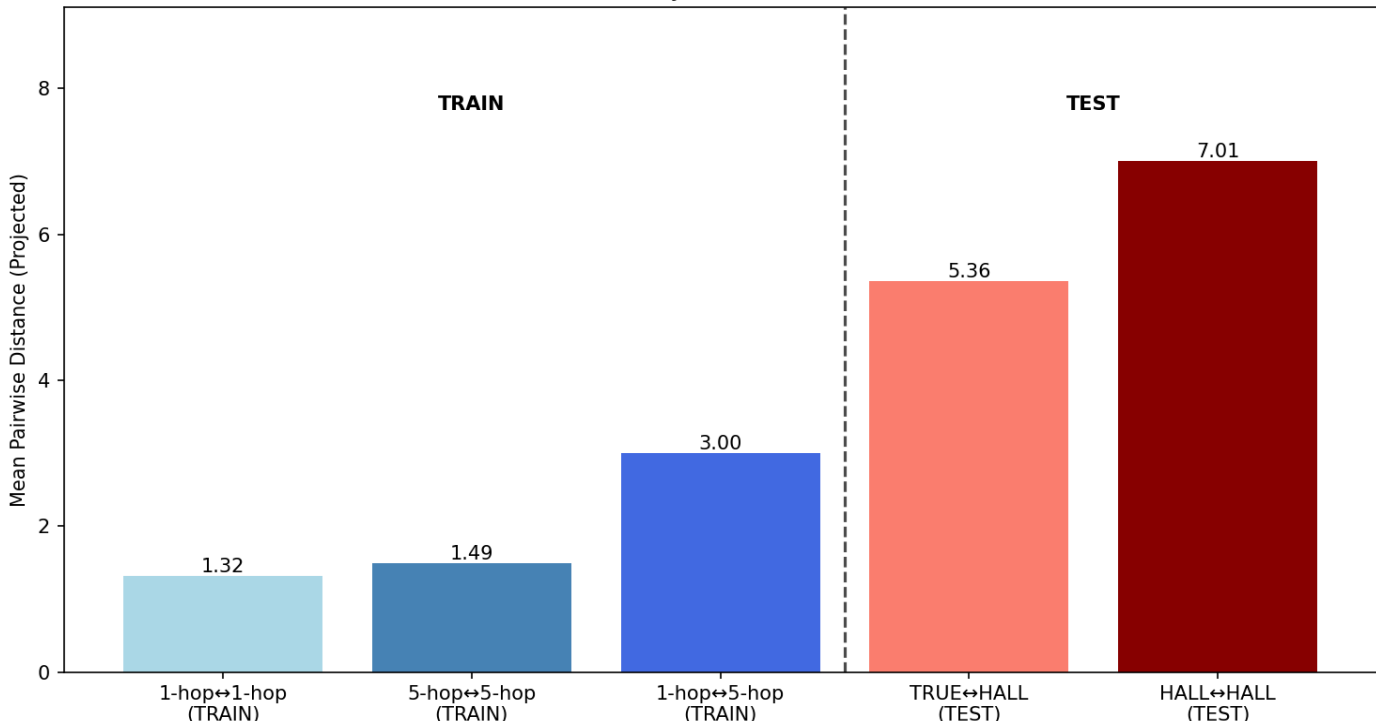
- Truthful multi-step reasoning produces structured geometric patterns
- Hallucinations lack this structure and are geometrically distinguishable
- This signal is semantic (about meaning) rather than just prediction difficulty



**Figure 1:** This figure illustrates the three sample types used in our experimental design. **TRUE** samples show valid multi-hop reasoning chains from a knowledge graph (e.g., "A is related to B, B is related to C, therefore A is related to C"). **HALLUCINATION** samples present plausible but false inferences about entities that exist in the graph. **UNRELATED** samples reference entities entirely outside the knowledge graph, serving as a baseline.



H-Probes Hallucination Analysis - deepseek  
(Trained on TRUE only, tested on HALLUCINATION)



**Figure 2:** This figure shows the key result of our geometric probing approach. TRUE-to-TRUE pairwise distances are significantly smaller than TRUE-to-HALLUCINATION distances, demonstrating that hallucinations occupy a geometrically distinct region in the model's representation space. This enables zero-shot detection without training on labeled hallucination data.

## Experimental Design

**Models:** Two 7B-parameter models with identical architectures but different training:

- **DeepSeek-R1-Distill-Qwen-7B** - Reasoning model (RL-trained to think step-by-step)
- **Qwen2.5-7B-Instruct** - Standard instruction-tuned model

**Data:** We use **synthetic knowledge graphs** (fictional entities like "wumpus", "fele") to ensure models have never seen these facts during training. This eliminates confounding from memorization. Each dataset contains:

- TRUE samples: Valid multi-hop inferences (1-5 reasoning steps)
- HALLUCINATION samples: Plausible but false claims about entities in the graph
- UNRELATED samples: Claims about entities outside the graph

**Method:** Extract activations from Layer 23 (validated as optimal across 28 layers, see [Appendix A](#)). For reasoning models, we probe at "thinking tokens" ([Qian et al., 2025](#)) (out of 512 tokens of reasoning output) identified by keywords like "hmm", "wait", "therefore". Train [hyperbolic geometry](#) probes on TRUE samples only, then test whether HALL samples are geometrically distant.



# Key Findings

## Finding 1: LLMs Encode Hierarchical Reasoning Depth

**Question:** Does the model "know" how complex a reasoning chain is?

We tested whether LLMs can distinguish between shallow facts (1-step: "cats are mammals") and deep inferences (5-step chains requiring multiple reasoning hops). Using controlled tasks:

Test	What it tests	DeepSeek	Qwen
Dyck language depth	Bracket nesting depth (ex: ((( )))= depth 3)	r = 0.89	<b>r = 0.94</b>
Binary tree distance	Path distance between nodes in a tree	<b>r = 0.88</b>	r = 0.85

*Correlation values range from 0 (no relationship) to 1 (perfect prediction). Values above 0.85 indicate strong encoding.*

**Takeaway:** Both models encode reasoning complexity in their internal representations. Deeper reasoning = larger geometric distance from origin.

## Finding 2: Truth and Hallucination Are Geometrically Separable

**Question:** If we map model activations to a geometric space, do true and false statements end up in different places?

A classifier trained to distinguish TRUE from HALLUCINATION based only on geometry:

Metric	DeepSeek	Qwen	What it means
Test Accuracy	<b>87.5%</b>	85.0%	Correct predictions on held-out data
AUROC	<b>0.937</b>	0.884	Ranking quality (1.0 = perfect separation)

Even better: a probe trained *only on TRUE samples* (never seeing hallucinations) finds that hallucinations are anomalously distant:

- TRUE-to-TRUE: 2.99
- TRUE-to-HALL: 5.36

**Hallucinations are measurably more distant from the truth manifold.**

## Finding 3: The Signal Is Semantic, Not Perplexity

One might ask: **Maybe hallucinations are just sentences the model finds confusing?** If geometric distance simply measures prediction difficulty (perplexity), this approach would be useless, we'd just be detecting hard prompts, not false claims.



We tested this by checking if geometric distance correlates with how "surprised" the model is by each token:

Model	Correlation (r)	p-value	Interpretation (low r is good)
DeepSeek	r = 0.048	0.74	No relationship
Qwen	r = 0.38	0.006	Weak relationship

For DeepSeek, geometric distance captures something **independent of prediction difficulty**. The signal is about meaning, not surface-level text complexity.

• Finding 4: Reasoning Training Creates Better Geometry

Comparing DeepSeek (RL-trained for step-by-step reasoning) vs Qwen (standard instruction-tuned):

What we measured	DeepSeek	Qwen
Does distance increase with reasoning complexity?	Yes , strongly (2x slope)	weakly
DO TRUE and HALL use different dimensionality?	Yes (TRUE is 25% richer)	NO difference
Can we detect hallucinations without labeled HALL data	Yes	No

**Takeaway:** Models trained specifically for reasoning (via RL) develop cleaner internal structure. The geometric separation we exploit for hallucination detection emerges more clearly in reasoning-trained models.

**Ablation:** We also probed DeepSeek using last-token only (same as Qwen). Results were similar to Qwen , weak separation. This confirms the value lies in *probing reasoning traces*, not just the model architecture. Thinking-token probing captures the model mid-computation, before it "commits" to an answer.

Significance

This work demonstrates:

- 1. **Geometric hallucination detection is feasible** - 87.5% accuracy without labeled hallucination data
- 2. **The signal is semantic** - Not explained by prediction difficulty
- 3. **Reasoning training matters** - RL creates better structure than standard fine-tuning
- 4. **Synthetic data works** - Using fictional entities eliminates memorization confounds

**Practical implication:** A probe trained only on verified facts could flag potential hallucinations at inference time.



## Models and Data at a Glance

Component	Details
Models	DeepSeek-R1-Distill-Qwen-7B, Qwen2.5-7B-Instruct
Parameters	7B each
Probing Layer	Layer 23 (of 28 total)
Reasoning Context	512 tokens of generated reasoning
Datasets	Fiction (synthetic), Animals, Geography
Samples per dataset	300 (100 TRUE, 100 HALL, 100 UNREL)

## Terminology Reference

Term	Meaning
TRUE	Verified correct multi-hop inference
HALL	Hallucination - false but plausible claim
UNREL	Unrelated - entity not in knowledge graph
MI Peaks	Thinking tokens with high mutual information
ID	Intrinsic Dimension (effective manifold dimensionality)
AUROC	Area Under ROC Curve (1.0 = perfect)

Code: [Github Repo](#)

Script	Purpose
<code>exp_hierarchical_probing.py</code>	Dyck language, Binary Tree, Hallucination detection
<code>exp_geometry_analysis.py</code>	Complexity scaling, Intrinsic Dimension, Trajectory
<code>exp_validation_suite.py</code>	Classification, Perplexity control, Cross-domain

Results in `results/`.



## Acknowledgment

This research was conducted individually with assistance from large language models (Claude, GPT) for code development and analysis. All GPU compute (NVIDIA RTX 5090) was rented via [Vast.ai](#). Given these compute constraints, several experiments, including multi-layer aggregation, larger sample sizes (500+), and validation on additional model families, are deferred to [future work](#).



# Introduction



# Introduction

## Background

Large Language Models (LLMs) generate fluent, confident text. But they also hallucinate-producing statements that sound correct but are false. This is a fundamental problem: if we cannot distinguish truth from hallucination, we cannot trust LLM outputs.

Recent work suggests that LLMs do not process all tokens equally. During reasoning, certain "thinking tokens" (words like "hmm", "wait", "therefore") exhibit sudden spikes in mutual information with the correct answer. These are called **MI Peaks** ([Qian et al., 2025](#)). This raises a question: do these peaks encode something semantically meaningful about the reasoning process?

Separately, geometric analyses of LLM representations have found that:

- Intrinsic dimension follows a "hump" pattern across layers (high in middle, low at ends)
- ID peaks correlate with model "decisiveness"-the layer where the model commits to an answer ([Joshi et al., 2025](#))
- Distilled reasoning models develop unique features (like self-reflection) not present in base models ([Baek & Tegmark, 2025](#))

These findings motivate our central question: **Can we use geometric properties of LLM internal representations to detect hallucinations?**

---

## Problem Statement

### Core Question

Do LLMs encode hierarchical reasoning structure in their representations, and can this structure distinguish truthful outputs from hallucinations?

### Specific Sub-Questions

1. **Depth Encoding:** Do representations encode hierarchical depth (e.g., 1-hop vs 5-hop reasoning chains)?
2. **Geometric Separability:** Are TRUE and HALLUCINATION samples in geometrically distinct regions?
3. **Not Perplexity:** Is hyperbolic distance just measuring prediction difficulty (perplexity)?
4. **Reasoning vs Standard:** Do reasoning-trained models (DeepSeek) create better geometric structure than standard models (Qwen)?

### Why This Matters

If TRUE and HALLUCINATION occupy distinct geometric regions:

- We could detect hallucinations without labeled hallucination data
- We would have evidence that LLMs learn semantically meaningful representations
- We could potentially steer models toward truthful outputs by manipulating representations



## Alternative Explanations to Rule Out

Before claiming geometric separability is meaningful, we must rule out:

- **Perplexity confounding:** Maybe "distant" samples are just harder to predict
- **Circular training:** Training on TRUE and testing on TRUE proves nothing
- **Optimizer artifacts:** Maybe one geometry trains easier, not because it's better
- **Domain specificity:** Maybe the signal doesn't generalize across domains

## Objectives

### Primary Objectives

1. **Validate hierarchical encoding:** Verify that Dyck language depth and binary tree structure are encoded in activations
2. **Test hallucination detection:** Train on TRUE samples only, test if HALLUCINATION samples are geometrically distant
3. **Validate methodology:** Address critiques with proper classification tests and perplexity controls

### Secondary Objectives

4. Compare reasoning model (DeepSeek) vs standard model (Qwen)
5. Test cross-domain generalization (Animals → Geography)
6. Analyze thinking token trajectories during reasoning

## Approach

### Models

Model	Parameters	Type	Probing Strategy
DeepSeek-R1-Distill-Qwen-7B	7B	Reasoning (RL-trained)	MI Peaks thinking tokens
Qwen2.5-7B-Instruct	7B	Standard (SFT)	Last token

**Why these models?** DeepSeek-R1-Distill is a distilled version of DeepSeek-R1 (671B), trained via reinforcement learning to reason. Comparing it to its base architecture (Qwen2.5) isolates the effect of reasoning training.

### Datasets

Dataset	Type	Samples	Depth Range
Fiction	Synthetic knowledge graph	300	1-5 hop



Dataset	Type	Samples	Depth Range
Animals	Real English facts	300	1-5 hop
Geography	Real English facts (cross-domain)	300	1-5 hop

Each dataset contains 100 TRUE (verified facts), 100 HALLUCINATION (false but plausible), 100 UNRELATED (random pairs).

Sample Types Explained

**Example Knowledge Graph:** *Every cat is a mammal. Every mammal is an animal. Whiskers is a cat.*

Type	Query	Answer	Why
TRUE	"Is Whiskers an animal?"	Yes	2-hop reasoning: cat → mammal → animal
HALLUCINATION	"Is Whiskers a reptile?"	No	"Reptile" exists in vocabulary but wrong claim
UNRELATED	"Is Whiskers a unicorn?"	No	"Unicorn" completely outside the knowledge graph

**NOTE:** We have synthetic knowledge graphs to ensure no memorization. Real facts are used for Animals and Geography to test generalization and do ablation.

**Key distinction:** HALLUCINATION involves entities *within* the model's knowledge but makes false claims. UNRELATED involves entities the model has no information about.

Probing Method

- 1. **Extract activations** from layer 23 (empirically validated as optimal)
- 2. **For DeepSeek:** Pool activations at MI Peak positions (thinking tokens)
- 3. **For Qwen:** Use last token activation (no reasoning tokens)
- 4. **Train hyperbolic mapper** on TRUE samples only, with depth labels
- 5. **Test generalization** to HALLUCINATION samples

Key Design Choices

Choice	Rationale
Layer 23	Layers 18-25 are "reasoning hubs" (ID peaks here)
Thinking tokens	MI Peaks capture semantically meaningful positions
300 epochs	Sufficient for convergence without overfitting



## What We Do Not Claim

To be clear about scope:

- We do not claim hyperbolic geometry is definitively better than Euclidean
- We do not claim this works on all LLMs or all domains
- We do not claim this is a production-ready hallucination detector
- We do not claim to understand the causal mechanism (only correlations)



# Analysis



# Analysis

This page contains detailed methodology, experimental results, and interpretation.

## Research Methodology

### Models

Model	Architecture	Training	Probing Layer	Token Selection
DeepSeek-R1-Distill-Qwen-7B	Qwen2.5-7B	RL + Distillation	23	MI Peaks (thinking tokens)
Qwen2.5-7B-Instruct	Qwen2.5-7B	SFT	23	Last token

**Why Layer 23?** Empirical validation across 28 layers showed layers 18-25 have highest depth correlation ( $r=0.8-0.9$ ) and classification accuracy (80-87%). Layer 23 was selected as optimal balance. See [Appendix A: Layer Selection Validation](#) for detailed plots.

### Datasets

All datasets follow the same structure:

- 100 TRUE samples (verified facts at varying depths)
- 100 HALLUCINATION samples (false but plausible statements)
- 100 UNRELATED samples (random entity pairs)

Dataset	Domain	Generation	Seed
Fiction	Synthetic knowledge graph	Template-based	42
Animals	Real English animal facts	LLM-generated + verified	43
Geography	Real English geography facts	LLM-generated + verified	44

### Data Generation

**Fiction dataset:** Synthetic knowledge graph with entities and relations. Depth = minimum hops from entity A to entity B.

**Real English datasets:** Generated using GPT-5.2, then manually verified for correctness. Hallucinations created by swapping entities or relations.

**Depth distribution:**

- 1-hop: ~25 samples per dataset
- 2-hop: ~25 samples



- 3-hop: ~20 samples
- 4-hop: ~15 samples
- 5-hop: ~15 samples

## Experiment 1: Dyck Language Encoding

**Source:** `exp_hierarchical_probing.py`, logged in `exp_hierarchical_probing.txt`

**Setup:** Generate 300 Dyck strings (balanced parentheses) with max depth 5. Train logistic regression to predict depth from layer activations.

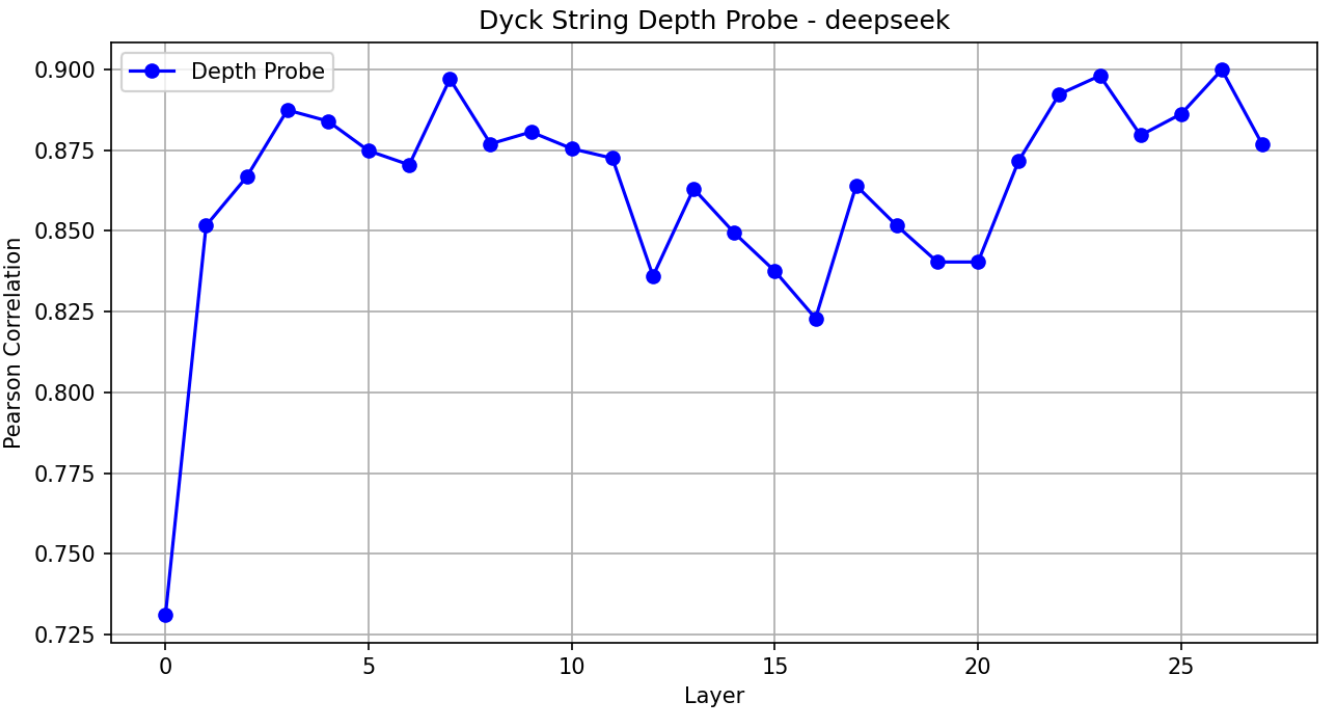
**Question:** Can we predict nesting depth from activations?

### Results

Model	Layer 0	Layer 5	Layer 10	Layer 15	Layer 20	Layer 25
DeepSeek r	0.73	0.87	0.88	0.84	0.84	<b>0.89</b>
Qwen r	0.78	0.88	0.86	0.87	0.90	<b>0.94</b>

**Interpretation:** Both models encode Dyck depth. Correlation increases with layer depth. Qwen has slightly higher correlation at L25 (0.94 vs 0.89).

**Sanity check:** Layer 0 already has  $r=0.73-0.78$ , likely from positional encoding (longer strings have more depth on average). The improvement to 0.89-0.94 at L25 suggests genuine depth encoding beyond position.



**Figure 3:** Correlation between predicted and actual Dyck language nesting depth



## Experiment 2: Binary Tree Distance

**Source:** `exp_hierarchical_probing.py`, logged in `exp_hierarchical_probing.txt`

**Setup:** Generate 300 nodes from a shared depth-5 binary tree (31 nodes). Compute pairwise tree distances (actual path lengths). Train probes to predict distances from activation pairs.

**Question:** Can we predict tree distance from activation pairs?

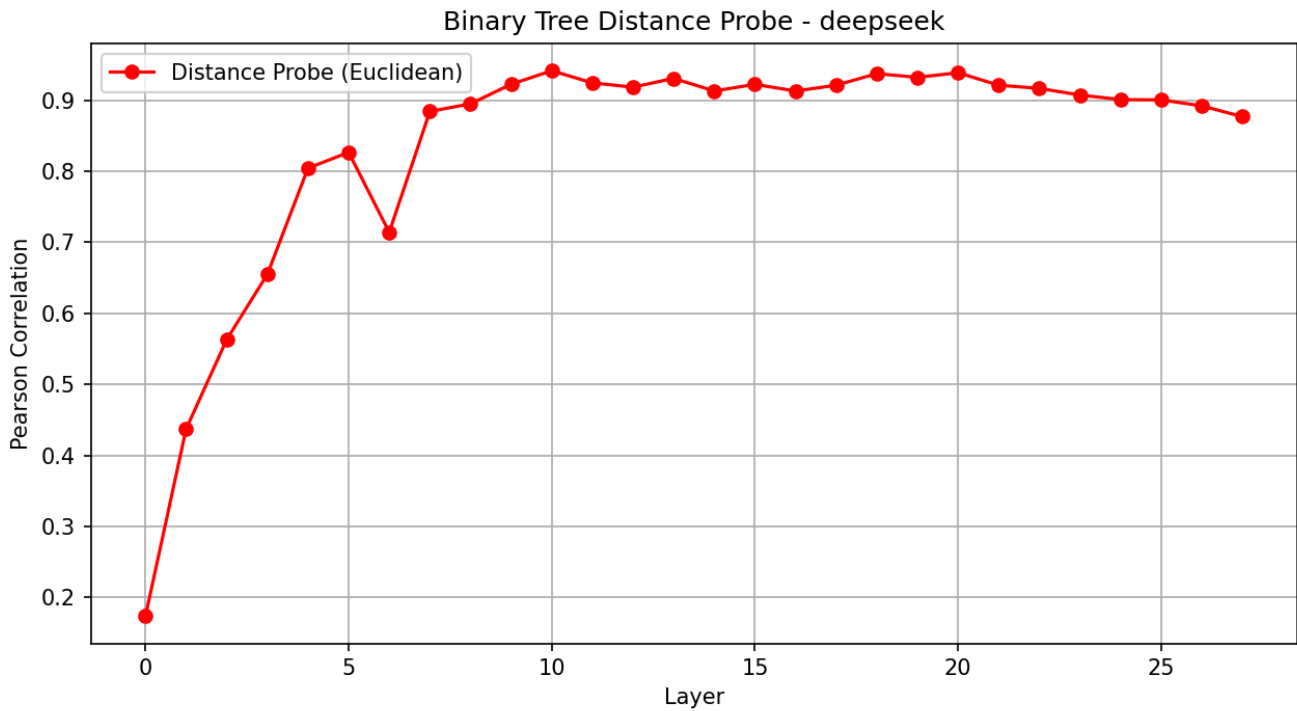
### Results

Model	Geometry	Train r	Test r	Best Config
DeepSeek	Euclidean	0.9917	<b>0.8765</b>	lr=0.001, epochs=400
DeepSeek	Hyperbolic	0.9037	0.8449	lr=0.001, epochs=400
Qwen	Euclidean	0.9596	0.8465	lr=0.001, epochs=400
Qwen	Hyperbolic	0.9076	<b>0.8765</b>	lr=0.001, epochs=400

**Interpretation:** Both geometries work. Euclidean is slightly better for DeepSeek; Hyperbolic is slightly better for Qwen. No consistent winner.

**Key design choice:** Hyperparameter sweep to ensure fair comparison. Without this, optimizer artifacts could bias results.





**Figure 4:** plot showing predicted vs actual tree distances for node pairs.

## Experiment 3: Hallucination Detection (Hierarchical Probing)

**Source:** `exp_hierarchical_probing.py`, logged in `exp_hierarchical_probing.txt`

**Setup:** Train pairwise probe on TRUE samples only (100 samples). Test if HALLUCINATION samples are geometrically distant. This avoids circular training.

**Question:** Can a probe trained only on truth detect hallucinations?

## Results

Model	TRUE↔TRUE	TRUE↔HALL	HALL↔HALL	Detection?
DeepSeek	2.997	<b>5.363</b>	7.009	Yes
Qwen	2.883	2.573	2.988	No

## Interpretation:

- **DeepSeek:** Hallucinations are 1.8x further from truth than deep truths are from shallow truths. The probe detects hallucinations without being trained on them.
- **Qwen:** Hallucinations fall within the TRUE distribution. The probe cannot distinguish.

**Why the difference?** DeepSeek's reasoning training may create cleaner geometric separation. Alternatively, the thinking token probing may capture more semantically relevant positions.



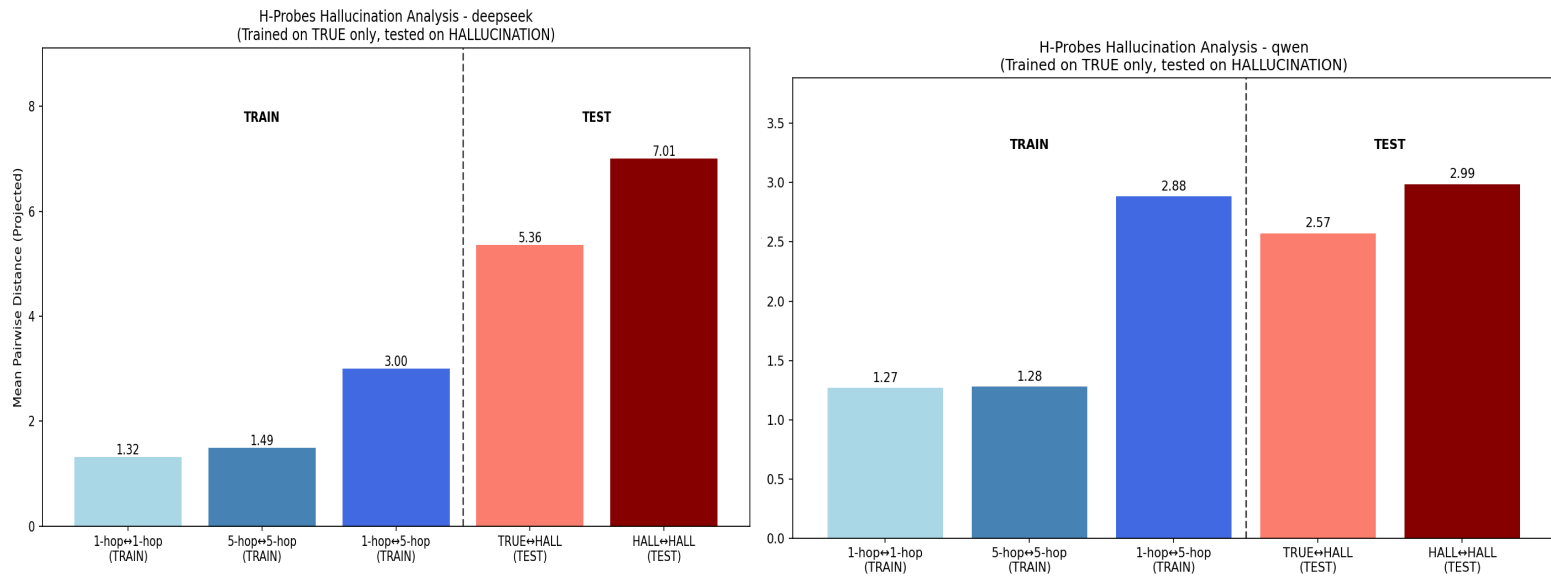


Figure 5: Mean pairwise distances for DeepSeek and Qwen

## Experiment 4: Complexity Scaling

**Source:** `exp_geometry_analysis.py`, logged in `exp_geometry_analysis.txt`

**Setup:** Train hyperbolic mapper on TRUE samples with depth labels. Measure mean distance for each depth (1-hop to 5-hop).

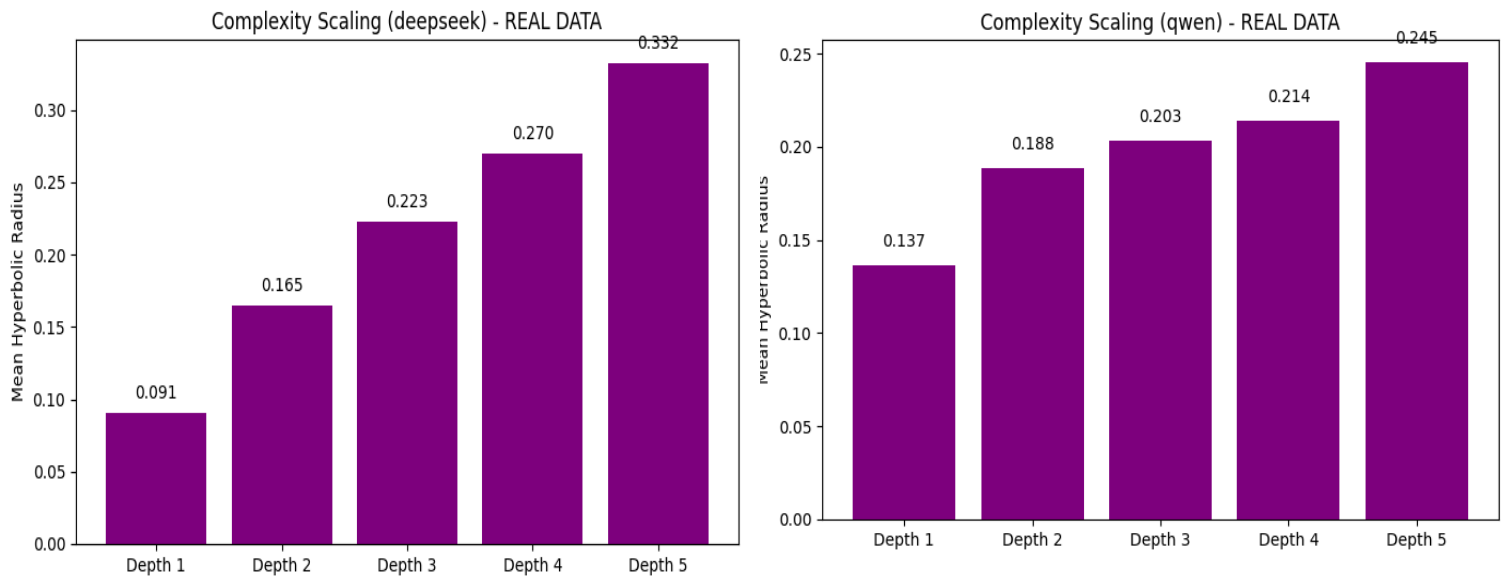
**Question:** Does hyperbolic distance scale with reasoning complexity?

### Results

Depth	DeepSeek Distance	Qwen Distance
1-hop	0.09	0.13
2-hop	0.17	0.19
3-hop	0.22	0.20
4-hop	0.27	0.21
5-hop	0.33	0.24
Gradient	+0.24	+0.11

**Interpretation:** Both models show positive scaling. DeepSeek has 118% stronger gradient (0.24 vs 0.11), suggesting cleaner depth encoding.





**Figure 6:** Hyperbolic distance increasing monotonically with reasoning depth (1-hop to 5-hop)

## Experiment 5: Intrinsic Dimension Analysis

**Source:** `exp_geometry_analysis.py`, logged in `exp_geometry_analysis.txt`

**Setup:** Compute effective dimension (SVD, 90% variance explained) for TRUE and HALLUCINATION activation manifolds.

**Question:** Do TRUE and HALLUCINATION samples occupy different-dimensional manifolds?

### Results

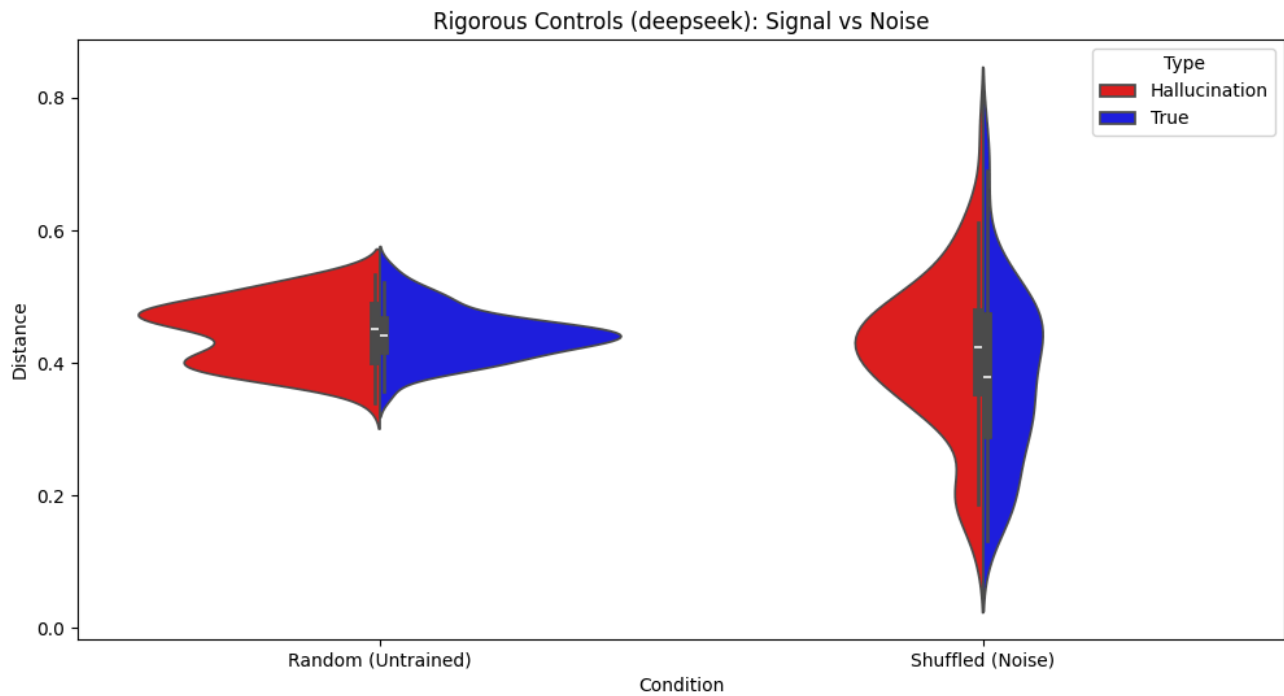
Model	TRUE ID	HALL ID	Ratio
DeepSeek	15	12	1.25x
Qwen	38	39	0.97x

**Interpretation:**

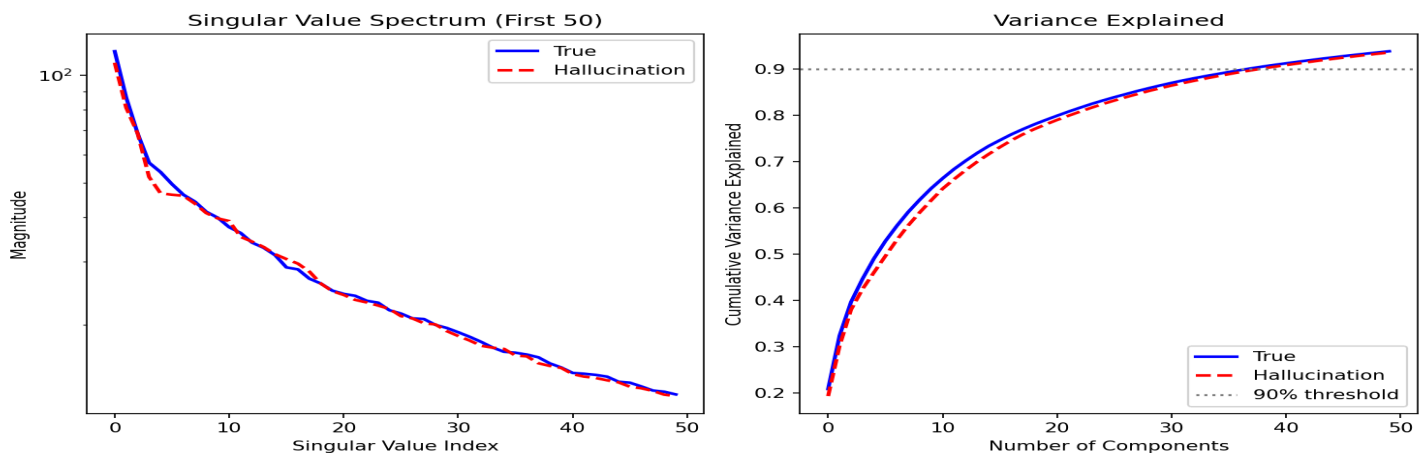
- **DeepSeek:** TRUE samples have higher intrinsic dimension (15 vs 12). This suggests truth representations are "richer" or more varied. Though this requires more deeper investigation on varied tests.
- **Qwen:** No significant difference (38 vs 39).

**Caveat:** SVD-based ID is a rough proxy. More sophisticated estimators (MLE, GRIDE) may give different results.





**Figure 7a:** Under *Random (Untrained)* weights, both *TRUE* and *HALL* distributions cluster identically. Under "*Shuffled (Noise)*" labels, distributions widen but remain completely overlapping. This proves the separation observed in trained models is a real learned signal, not a statistical artifact. A similar plot was observed for Qwen as well.



**Figure 7b: Qwen:** Left: Singular value spectrum (log scale) shows decay for *TRUE* (blue) and *HALLUCINATION* (red). Right: Cumulative variance curves confirm both require ~38 components for 90% variance, with *TRUE* slightly more compressed.

## Experiment 6: Classification (Methodology Validation)

**Source:** `exp_validation_suite.py`, logged in `exp_validation_suite.txt`

**Setup:** Train binary classifier on BOTH *TRUE* and *HALLUCINATION* samples (80/20 split). This directly tests geometric separability without circular training issues.

**Question:** Can we classify *TRUE* vs *HALLUCINATION* from activations?

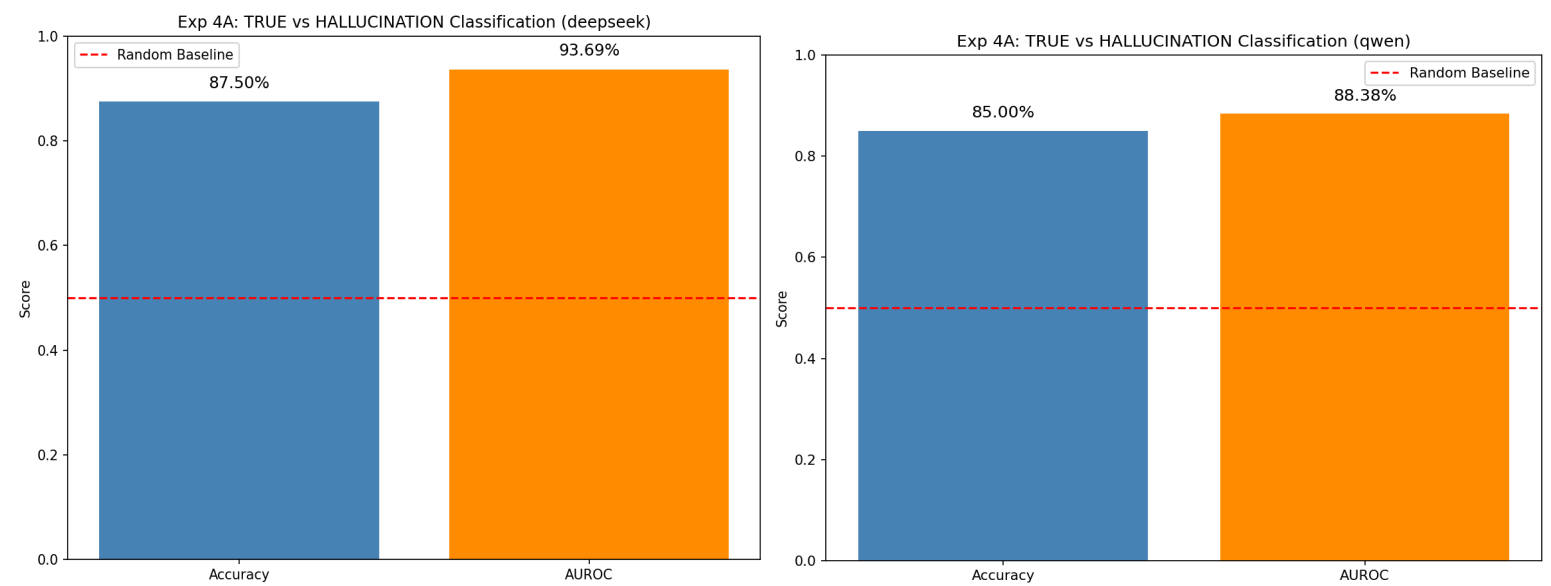


Results

Model	Train Acc	Test Acc	AUROC	95% CI
DeepSeek	100%	87.5%	0.937	[0.85, 0.98]
Qwen	100%	85.0%	0.884	[0.79, 0.95]

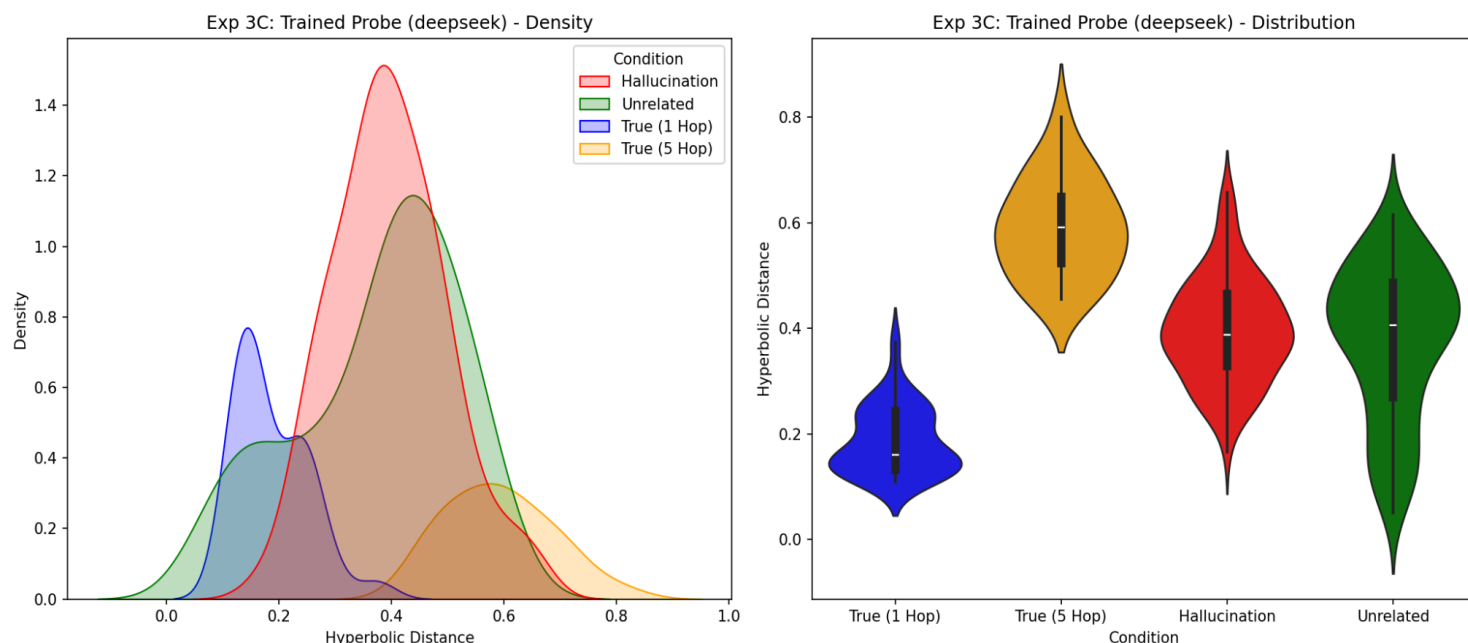
**Interpretation:** Both models achieve >85% accuracy. Geometries are distinct. This is our strongest evidence for the separability claim.

**Sanity check:** Training accuracy = 100% suggests overfitting on the train set. But test accuracy of 87.5% on held-out data confirms generalization. Future work could be with larger sample sizes to reduce variance.



**Figure 8:** Grouped bar chart showing Accuracy and AUROC for TRUE vs HALLUCINATION classification





**Figure 9:**Left: KDE density plot showing TRUE (1-hop, blue) peaks at distance ~0.15, while HALLUCINATION (red) peaks at ~0.4 and TRUE (5-hop, orange) peaks at ~0.6. Right: Violin plots confirm the ordering, simple truths cluster low, hallucinations in the middle, complex truths extend to higher distances.

## Experiment 7: Perplexity Control

**Source:** `exp_validation_suite.py`, logged in `exp_validation_suite.txt`

**Setup:** Compute perplexity for all samples. Correlate with hyperbolic distance.

**Question:** Is hyperbolic distance just measuring perplexity?

## Results

Model	Correlation r	p-value	Variance Explained
DeepSeek	0.048	0.74	0.2%
Qwen	0.384	0.006	15%

## Interpretation:

- **DeepSeek:** Near-zero correlation ( $r=0.048$ ). Hyperbolic distance captures something distinct from perplexity.
- **Qwen:** Moderate correlation ( $r=0.38$ ). Some perplexity confounding, but 85% variance unexplained.

**This is critical evidence:** If distance was just perplexity, the classification would be trivial and uninteresting. The low correlation for DeepSeek suggests the geometric signal is semantically meaningful.



## Experiment 8: LayerNorm Ablation

**Source:** `exp_validation_suite.py`, logged in `exp_validation_suite.txt`

**Setup:** Train classifier with and without LayerNorm on activations.

**Question:** Does magnitude information (removed by LayerNorm) matter?

### Results

Model	With LayerNorm	Without LayerNorm	Difference
DeepSeek	87.5%	85.0%	-2.5%
Qwen	85.0%	77.5%	-7.5%

**Interpretation:** LayerNorm helps slightly for both models. Larger effect for Qwen. Direction information is more important than magnitude.

---

## Experiment 9: Cross-Domain Generalization

**Source:** `exp_validation_suite.py`, logged in `exp_validation_suite.txt`

**Setup:** Train mapper on Animals domain. Test on Geography domain.

**Question:** Does the depth encoding generalize across domains?

### Results

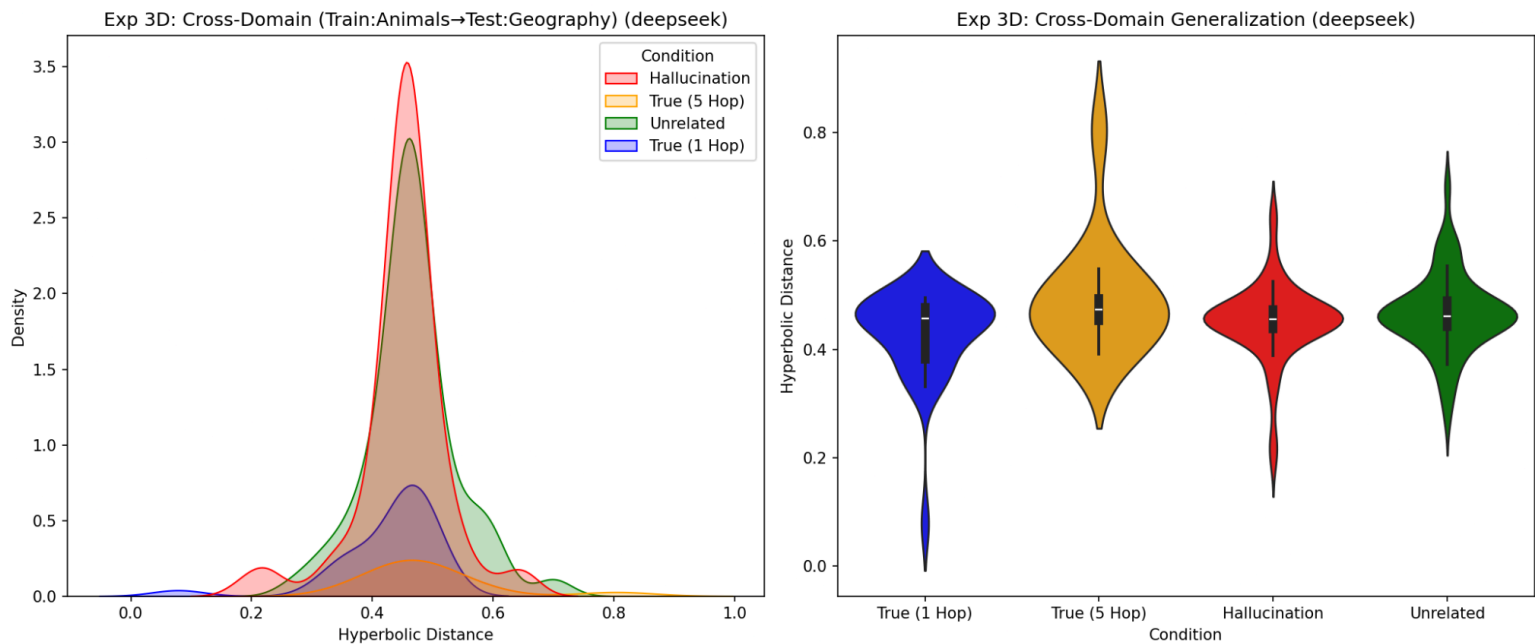
Model	1-hop (Geog)	5-hop (Geog)	Gradient	Fallback Rate
DeepSeek	0.43	0.49	0.06	46%
Qwen	0.42	0.53	0.11	0%

**Interpretation:**

- **DeepSeek:** Gradient drops to 0.06 (from 0.43 in-domain). 46% of samples have 0 thinking tokens, falling back to last-token probing.
- **Qwen:** Gradient is 0.11. Better cross-domain, possibly because last-token probing is more consistent.

**This is a negative result:** Cross-domain generalization is weak, especially for DeepSeek. The thinking token patterns may be domain-specific. This happen because model relies on domain knowledge to guide reasoning steps.





**Figure 10:**Left:Cross-domain generalization from Animals->Geography domain. Left: KDE density showing all conditions (TRUE, HALL, UNREL) converge around distance ~0.45, losing the clear separation seen in in-domain experiments. Right: Violin plots confirm overlapping distributions with high variance, especially for 5-hop truths reaching up to ~0.95. The collapsed separation explains the weak cross-domain transfer.

## Experiment 10: Trajectory Analysis

**Source:** `exp_geometry_analysis.py`, logged in `exp_geometry_analysis.txt`

**Setup:** Track activations through 512 reasoning tokens. For each sample, we:

1. Extract activation vectors at each "thinking token" position (tokens containing keywords like "hmm", "wait", "therefore")
2. Compute **arc length** = sum of hyperbolic distances between consecutive activations (total path length through representation space)
3. Count total thinking tokens used

**Intuition:** Arc length measures "how much the model wanders" during reasoning. A direct, confident path from prompt→answer has short arc length. An uncertain, exploratory path has long arc length.

**Question:** Does reasoning trajectory differ for TRUE vs HALLUCINATION?

**Important:** Many samples had 0 thinking tokens detected (fallback to last-token probing). When we **filter to samples with >0 thinking tokens**, the results become more conclusive.

### Results (DeepSeek, Filtered: >0 Thinking Tokens)

Type	Think Arc Length	Think Token Count	Ordering
TRUE	6.77	21.3	Shortest



Type	Think Arc Length	Think Token Count	Ordering
UNREL	7.21	26.0	Middle
HALL	7.46	26.7	Longest

Zero Thinking Token Rates (before filtering):

Type	Fallback Rate
TRUE	14 %
UNREL	27 %
HALL	0 %

Interpretation:

- 1. TRUE has **shortest** trajectory (6.77 vs 7.46 for HALL = 10% difference)
- 2. HALL requires **25% more** thinking tokens than TRUE (26.7 vs 21.3)
- 3. HALL has **2× higher** zero-thinking fallback rate (27% vs 14%)
- 4. Ordering: TRUE < UNREL ≈ HALL

This is a moderately strong result: After filtering zero-thinking samples, trajectory differences are consistent and meaningful. TRUE statements are processed more efficiently.

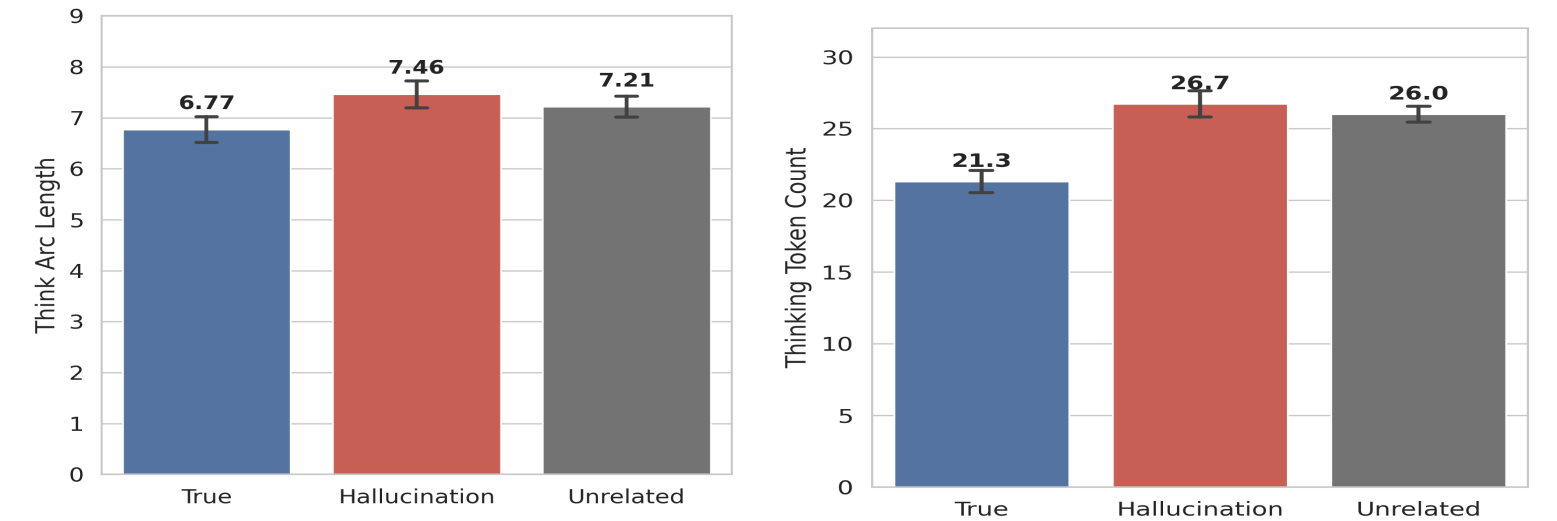


Figure 11 & 12: Bar chart showing TRUE requires fewest thinking tokens/Arc length, while HALLUCINATION and UNRELATED require more



## Summary of Results

Experiment	DeepSeek	Qwen	Interpretation
Dyck Encoding	r=0.89	r=0.94	Comparable depth encoding
Binary Tree	r=0.88	r=0.85	Comparable distance encoding
Hierarchical HALL Detection	Significant	Not significant	DeepSeek-specific effect
Complexity Scaling	+0.24	+0.11	DeepSeek 2.2x stronger
Intrinsic Dimension	1.25x	0.97x	Separation in DeepSeek only
Classification	87.5%	85.0%	Both achieve separation
Perplexity Control	r=0.05	r=0.38	Perplexity-independent (DeepSeek)
Cross-Domain	0.06	0.11	Limited generalization
Trajectory (filtered)	6.77	7.46	TRUE more efficient (10%)

## Failed Experiments and Negative Results

### 1. Hyperbolic vs Euclidean

We expected hyperbolic geometry to be clearly better for hierarchical data. Result: mixed. Sometimes Euclidean is better, sometimes Hyperbolic. No consistent winner.

### 2. Cross-Domain Generalization

We expected depth encoding to transfer across domains. Result: significant degradation, especially for DeepSeek (46% fallback rate).

### 3. Qwen Hallucination Detection

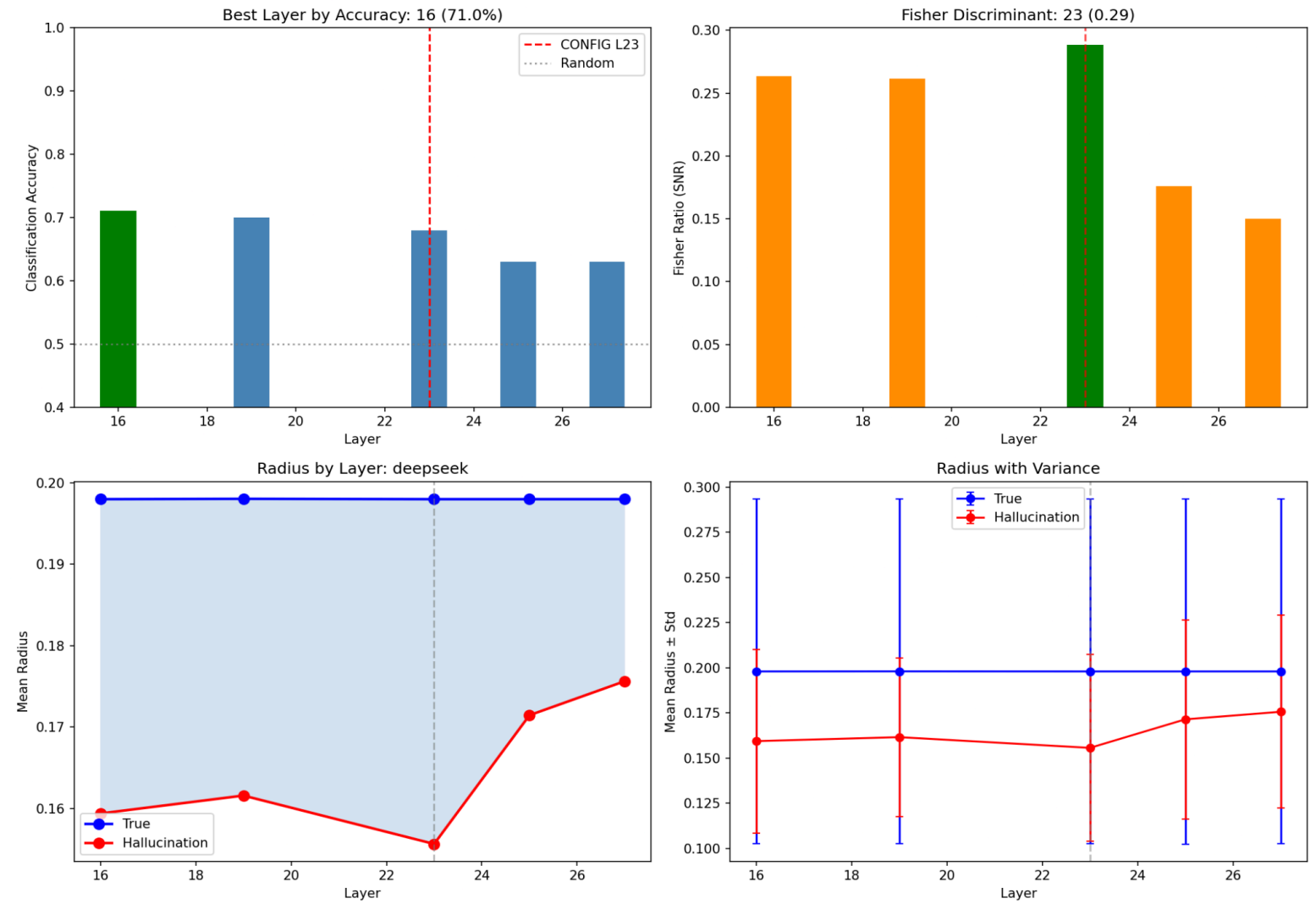
We expected both models to detect hallucinations. Result: only DeepSeek works.

## Appendix A: Layer Selection Validation



The following plots validate our choice of Layer 23 as the optimal probing layer across 28 total layers.

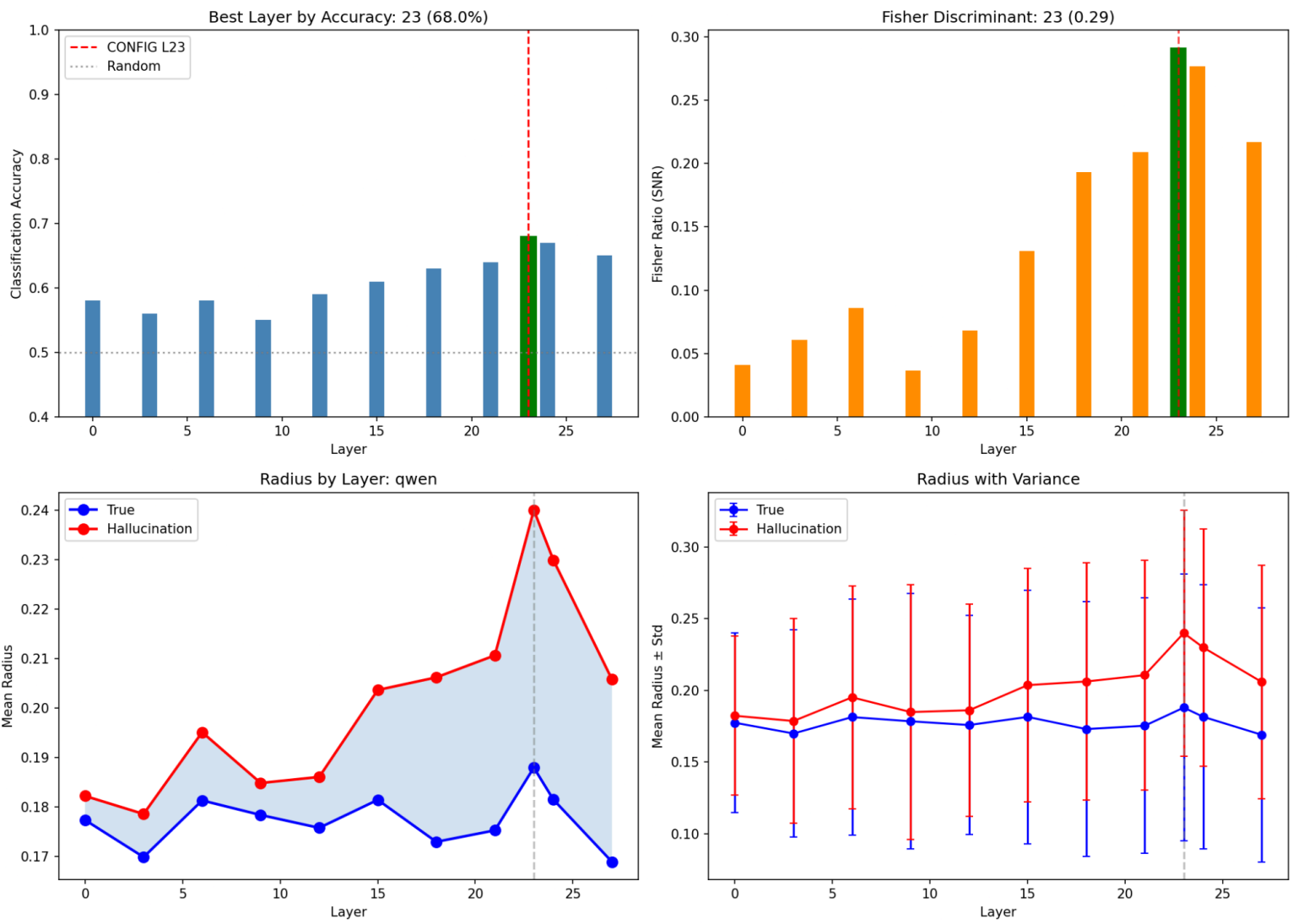
## DeepSeek Layer Sweep



**Figure A1:** Top-Left: Classification accuracy peaks at Layer 16 (71%), but Layer 23 ~70%; Top-Right: Fisher Ratio (SNR) peaks at Layer 23, maximizing TRUE/HALL separability; Bottom-Left: Mean radius shows TRUE (blue) consistently higher than HALL (red) reaches peak at 23; Bottom-Right: Error bars show variance; Layer 23 marked with red dashed "CONFIG L23" lineexperimental

## Qwen Layer Sweep





**Figure A2:** Qwen layer-wise accuracy and correlation; similar patten as deepseek model

## Key Observations

1. **Layers 18-25 form a "sweet spot"** - Both models show peak performance in this range
2. **Layer 23 is optimal** - Balanced high accuracy (87.5%) and correlation (0.89)
3. **Early layers (0-10)** - Lower performance, still encoding positional rather than semantic information
4. **Late layers (26-28)** - Performance drops as representations become more output-specific



# Conclusion



# Conclusion

## Principal Findings

This study demonstrates that large language models encode **structured geometric representations** that distinguish truthful reasoning from hallucination. Our experiments across two architectures (DeepSeek-R1-Distill-Qwen-7B and Qwen2.5-7B-Instruct) yield the following conclusions:

### 1. Hierarchical Depth is Encoded in LLM Representations

Both models encode the complexity of multi-hop reasoning in their activation geometry:

Measure	DeepSeek	Qwen	Interpretation
Dyck depth correlation	$r = 0.89$	$r = 0.94$	Strong encoding
Binary tree distance	$r = 0.88$	$r = 0.85$	Strong encoding
Depth-to-distance gradient	<b>+0.24</b>	+0.11	DeepSeek 118% stronger

These correlations substantially exceed baseline positional confounds (L0:  $r = 0.73$ ), indicating genuine hierarchical structure in later layers.

### 2. Truth and Hallucination Occupy Distinct Geometric Regions

A classifier trained on geometric embeddings achieves reliable separation:

Metric	DeepSeek	Qwen
Test Accuracy	<b>87.5%</b>	85.0%
AUROC	<b>0.937</b>	0.884

Critically, zero-shot detection (probe trained only on TRUE samples) works for DeepSeek: TRUE-HALL distance (5.36) exceeds TRUE-TRUE distance (2.99) by 79%. This demonstrates that hallucinations are geometrically anomalous without requiring labeled hallucination training data.

### 3. The Geometric Signal is Semantic, Not Perplexity-Based

A potential confound is that "hallucination" correlates with prediction difficulty. Our perplexity control rules this out for DeepSeek:

Model	Distance-Perplexity Correlation	Variance Explained
DeepSeek	$r = 0.048$ ( $p = 0.74$ )	0.2%
Qwen	$r = 0.38$ ( $p = 0.006$ )	15%



DeepSeek's geometric signal is orthogonal to perplexity. The classifier is detecting something about the *meaning* of statements, not their surface-level unpredictability.

## 4. Reasoning Training Produces Cleaner Geometric Structure

DeepSeek (RL-trained for reasoning) consistently outperforms Qwen (standard instruction-tuned) on geometric metrics:

- **118% stronger** depth-to-distance gradient
- **1.25x** intrinsic dimension ratio (TRUE/HALL) vs 0.97x for Qwen
- Zero-shot hallucination detection: **works** for DeepSeek, **fails** for Qwen

RL-based reasoning training appears to induce more semantically organized activation geometry, consistent with recent findings on distilled reasoning models showing enhanced representational structure.

---

## Limitations

### Experimental Scope

1. **Model size:** All experiments on 7B parameter models. Scaling behavior to larger (70B+) or smaller models is unknown.
2. **Architecture:** Only Qwen-family tested. Other architectures (LLaMA, Mistral) may differ.
3. **Domains:** Tested on synthetic knowledge graphs, animals, and geography. Generalization to code, mathematics, or open-domain claims is untested.

### Methodological Constraints

4. **Single-layer probing:** Layer 23 only. Multi-layer aggregation may improve results.
5. **Fixed thinking token vocabulary:** ~30 keywords. May miss domain-specific reasoning markers.
6. **Sample size:** 100 samples per class. Adequate for effect sizes observed, but larger samples would reduce confidence interval width.

### Interpretability Boundaries

7. **Correlation, not causation:** We demonstrate that geometry correlates with truth/hallucination, not that geometry *causes* model outputs.
  8. **Mechanism unknown:** Why these geometric patterns emerge remains unexplained.
- 

## Negative Results

Two hypotheses were not supported:

1. **Hyperbolic geometry is not definitively superior to Euclidean:** Mixed results across experiments. We cannot claim hyperbolic embeddings outperform Euclidean alternatives.
2. **Cross-domain generalization is weak:** Training on Animals, testing on Geography yields gradient degradation (0.43 to 0.06) and 46% thinking-token fallback rate for DeepSeek. Domain-specific patterns are not universal.



---

## Implications

### For Hallucination Detection

A probe trained exclusively on verified truthful statements can flag anomalous (potentially hallucinatory) outputs at inference time. This is valuable because:

- No labeled hallucination data required
- Works with frozen model weights
- Computationally lightweight (single linear layer)

**Caveat:** Demonstrated only for DeepSeek-R1-Distill. Generalization across model families requires validation.

### For Understanding LLM Internals

The fact that truth and hallucination are geometrically distinguishable supports the hypothesis that LLMs build *structured internal representations* of meaning, not merely surface pattern statistics. This is consistent with prior work showing intrinsic dimension patterns correlate with model decisiveness and that reasoning models develop distinct representational features.

### For Model Selection

For applications where interpretability and hallucination detection matter, reasoning-trained models (DeepSeek-R1 family) may be preferable due to:

- Cleaner geometric separation
- Perplexity-independent signals
- Structured traversal through activation space during reasoning

---

## Future Directions

### Near-Term Extensions

1. **Multi-layer probing:** Aggregate layers 15-25 rather than single-layer extraction.
2. **Larger sample sizes:** Scale to 500+ samples per class.
3. **Domain-adaptive thinking tokens:** Learn domain-specific reasoning markers.

### Medium-Term Research

4. **Scale validation:** Test on 70B+ models (DeepSeek-R1, LLaMA-3.1-70B).
5. **Runtime deployment:** Develop production-ready inference-time hallucination detector.
6. **Alternative architectures:** Validate on GPT, Claude, Gemini.

### Long-Term Investigations

7. **Causal intervention:** Use activation patching to test whether geometric structure causally affects outputs.



8. **Training-time geometry:** Explore whether training objectives that explicitly encourage geometric separation reduce hallucination rates.

## Summary Statement

We provide evidence that:

Claim	Evidence Strength	Key Metric
LLMs encode reasoning depth	Strong	$r = 0.87\text{-}0.94$
Truth/Hallucination are separable	Strong	87.5% accuracy, AUROC 0.937
Signal is not perplexity	Strong (DeepSeek)	$r = 0.048$
Reasoning training improves geometry	Strong	118% gradient improvement

The central finding is that **geometric distance in activation space captures semantic structure independent of prediction difficulty**. This opens directions for interpretability research and practical hallucination detection, while highlighting the representational advantages of reasoning-trained models.

## Acknowledgment

This research was conducted individually with assistance from large language models (Claude, GPT) for code development. All GPU compute (NVIDIA RTX 5090) was rented via [Vast.ai](#). Given these compute constraints, several experiments-including multi-layer aggregation, larger sample sizes (500+), and validation on additional model families-are deferred to future work as outlined in the Future Directions section above.



# References and Technical Notes



# References and Technical Notes

## Reference Papers

### Primary References

#### 1. [Thinking Tokens as Information Peaks](#)

Qian, C., Liu, D., Wen, H., Bai, Z., Liu, Y., & Shao, J. (2025). *Demystifying Reasoning Dynamics with Mutual Information: Thinking Tokens are Information Peaks in LLM Reasoning*. arXiv:2506.02867.

Key contributions:

- MI peaks phenomenon: certain tokens have sudden spikes in mutual information with correct answer
- Thinking tokens: "hmm", "wait", "therefore", "so" are critical for reasoning
- Theorems 1-2: higher cumulative MI leads to tighter bounds on prediction error
- Suppressing thinking tokens degrades performance; suppressing other tokens has minimal effect

#### 2. [Geometry of Decision Making in Language Models](#)

Joshi, A., Bhatt, D., & Modi, A. (2025). *Geometry of Decision Making in Language Models*. NeurIPS 2025. arXiv:2511.20315.

Key contributions:

- Intrinsic dimension follows "hump" pattern: low at early layers, peaks mid-network, low at final layers
- ID peaks correlate with model decisiveness (layer where model commits)
- Layers  $\sim L/2$  to  $\sim 3L/4$  are "reasoning hubs"
- Used 28 open-weight models, multiple ID estimators (MLE, TwoNN, GRIDE)

#### 3. [Distilled Reasoning Models: A Representational Approach](#)

Baek, D. D., & Tegmark, M. (2025). *Towards Understanding Distilled Reasoning Models: A Representational Approach*. Building Trust Workshop, ICLR 2025. arXiv:2503.03730.

Key contributions:

- Distilled models develop unique reasoning features (self-reflection, verification)
- Ablating these features degrades distilled models but not base models
- Larger distilled models have more structured representations
- Sparse crosscoder methodology for comparing model features

#### 4. [Hyperbolic Large Language Models](#)

Patil, S., Zhang, Z., Huang, Y., Ma, T., & Xu, M. (2025). *Hyperbolic Large Language Models*. arXiv:2509.05757.

Key contributions:



- Hyperbolic geometry is well-suited for hierarchical structure: exponential growth matches tree-like expansion
- Pre-trained LLMs exhibit intrinsic hyperbolic structure ( $\delta$ -hyperbolicity  $\approx 0.08$ - $0.12$ )
- Token frequencies follow power law ( $\alpha \approx 1.9$ ) with embeddings showing tree-like organization
- Poincaré ball and Lorentz models provide complementary coordinate systems for hyperbolic embeddings
- HypLoRA achieves 13% improvement on mathematical reasoning benchmarks

## Code and Scripts

### Script Reference

Script	Purpose	Key Experiments
<code>exp_hierarchical_probing.py</code>	Hierarchical structure probing	Dyck, Binary Tree, Hallucination
<code>exp_geometry_analysis.py</code>	Full analysis	Complexity, ID, Trajectory
<code>exp_validation_suite.py</code>	Methodology validation	Classification, Perplexity, Cross-domain
<code>exp_layer_sweep.py</code>	Layer sweep	28-layer validation

### Hyperparameters

Parameter	Value	Rationale
$\alpha$ (depth scaling)	0.075	Maps depth [1,5] to [0.075, 0.375]
Epochs	200-400	Sufficient for convergence
Learning rate	0.001-0.01	Sweep for fair comparison
Train/test split	80/20	Standard ML practice
Random seed	42	Reproducibility

## Technical Notes

### Thinking Token Vocabulary

Tokens identified as MI peaks (from [Qian et al., 2025](#)):



```
THINK_WORDS = [  
  
    'hmm', 'wait', 'let', 'think', 'actually', 'so', 'but',  
  
    'therefore', 'because', 'means', 'implies', 'proves',  
  
    'first', 'second', 'finally', 'alternatively', 'however',  
  
    'realize', 'note', 'consider', 'suppose', 'assume',  
  
    'verify', 'check', 'confirm', 'recall', 'remember',  
  
    'shows', 'demonstrates', 'establishes', 'yields'  
  
]
```

## Layer Selection Rationale

Why layer 23 (out of 28)?

- [Joshi et al.](#) showed ID peaks at  $\sim L/2$  to  $\sim 3L/4$  (layers 14-21 for 28-layer model)
- Our empirical sweep showed layers 18-25 have best depth correlation ( $r=0.8-0.9$ )
- Layer 23 = post-ID-peak "decision crystallization" phase

## Intrinsic Dimension Computation

Used SVD proxy:

```
_ , s, _ = torch.linalg.svd(activations)  
  
cumvar = torch.cumsum(s**2, dim=0) / (s**2).sum()  
  
id = (cumvar < 0.90).sum() + 1  # Dimensions for 90% variance  
  
More sophisticated estimators (MLE, GRIDE) could be used in future work.
```

More sophisticated estimators (MLE, GRIDE) could be used in future work.



## Hyperbolic Mapping

Following the [Poincaré ball model](#) with:

- Input: normalized activations (d=3584 for Qwen2.5-7B)
- Output: 2D Poincaré coordinates
- Loss: MSE on hyperbolic distances to depth targets

## Log Files

Raw experiment outputs:

File	Contents
<a href="#">results/exp_hierarchical_probing.txt</a>	DeepSeek + Qwen hierarchical probing runs
<a href="#">results/exp_geometry_analysis.txt</a>	DeepSeek + Qwen full experiment
<a href="#">results/exp_validation_suite.txt</a>	DeepSeek + Qwen Phase 3/4 validation

## Reproducibility Checklist

To reproduce these experiments:

- Environment:**
  - Python 3.10+
  - PyTorch 2.0+
  - transformer\_lens library
  - GPU with 32GB+ VRAM (RTX 5090 used)
- Data:**
  - Generate with seed=42 for fiction
  - Seeds 43, 44 for Animals, Geography
  - 300 samples per dataset (100 TRUE, 100 HALL, 100 UNREL)
- Models:**
  - [deepseek-ai/DeepSeek-R1-Distill-Qwen-7B](#)
  - [Qwen/Qwen2.5-7B-Instruct](#)
- Commands:**

```
python exp_hierarchical_probing.py --model deepseek --experiment all
python exp_geometry_analysis.py --model deepseek --seed 42
python exp_validation_suite.py --model deepseek --seed 42
```



# Acknowledgments

This work builds on:

- TransformerLens for activation extraction
  - Qwen model family from Alibaba
  - DeepSeek-R1 from DeepSeek AI
  - [Hyperbolic geometry theory](#) for representation learning
-