# Ensemble Verification for LLM Output Quality Assessment: **Lessons from the Synthetic-to-Production Gap**

Roman Khokhla
Independent Researcher
`rkhokhla@gmail.com`

October 25, 2025

**Abstract**

The discovery that compressibility-based signals achieve perfect detection (AUROC 1.000) on synthetic degeneracy but flag high-quality outputs on production models (GPT-4) reveals a fundamental challenge: **different failure modes require different signals**. We investigate ensemble approaches combining geometric signals ($r_{\text{LZ}}$ compressibility, lexical diversity, sentence repetition) with perplexity-based methods for comprehensive quality assessment.

Through analysis of 8,071 real GPT-4 outputs from production benchmarks (TruthfulQA, FEVER, HaluEval), we find: (1) Signal complementarity: perplexity captures factual errors; geometric signals capture structural patterns; lexical diversity correlates with sophistication. (2) Production reality: modern LLMs (GPT-4) avoid synthetic benchmark failures; signals designed for degraded models don't transfer. (3) Ensemble limitations: without ground-truth hallucination labels distinguishing factual errors from structural issues, ensemble methods cannot be rigorously validated.

This paper presents methodology, findings, and honest limitations, emphasizing the need for multi-modal evaluation frameworks that account for model evolution.

# 1 Motivation: Why Ensemble Approaches?

## 1.1 The Multi-Modal Nature of LLM Failures

LLM outputs can fail in fundamentally different ways:

- **Factual errors**: Incorrect claims, false information, contradicting known facts

- **Structural pathology**: Repetitive loops, semantic drift, incoherence

- **Quality degradation**: Poor lexical variety, simplistic language, hedging

Each failure mode has distinct signatures requiring specialized detection:

- **Factual errors** → Perplexity, NLI entailment, retrieval-augmented verification

- **Structural pathology** → Compression ratio ($r_{\text{LZ}}$), repetition detection

- **Quality markers** → Lexical diversity, coherence metrics

## 1.2 The Synthetic-Production Gap Challenge

Our previous work [1] discovered that:

- Compressibility signal ($r_{LZ}$) achieves **AUROC 1.000** on synthetic degeneracy

- Same signal on 8,290 real GPT-4 outputs flags **high-quality** responses (inverse enrichment)

- Outliers exhibit **higher** lexical diversity (0.932 vs 0.842, Cohen's $d = 0.90$)

- Outliers exhibit **lower** sentence repetition (0.183 vs 0.274, Cohen's $d = -0.47$)

**Interpretation**: Modern production models (GPT-4) are trained so well they don't produce the structural pathologies that synthetic benchmarks assume. Geometric signals detect what compresses—but in production, **sophistication** compresses as efficiently as **degeneracy** (for opposite reasons).

## 1.3 Research Questions

Given these findings, we investigate:

1. Can ensemble methods combining perplexity + geometric signals outperform perplexity alone?

2. Do different signals correlate with different failure modes in production outputs?

3. What are the limitations of ensemble approaches when models avoid synthetic failures?

## 2 Related Work

**Perplexity-based detection**: Simple, fast, proven for factuality [2]. AUROC $\sim$0.615 on factual hallucinations. Fails on structural degeneracy (AUROC 0.018, inverse correlation with confidence).

**Geometric/statistical methods**: SelfCheckGPT [4]: Sample consistency via NLI. $r_{LZ}$ compressibility: Perfect on synthetic, limited utility on GPT-4 (our work). Lexical diversity: Correlates with quality, not pathology.

**Ensemble approaches**: G-Eval [5]: GPT-4-as-judge with chain-of-thought. Multi-signal voting: Combines diverse signals but requires labeled data. Challenge: No public benchmarks with fine-grained failure mode labels.

## 3 Methodology

### 3.1 Data

**8,071 real GPT-4 outputs** (filtered, $n \geq 10$ tokens) from:

- **TruthfulQA** (790 samples): Misconceptions, false beliefs

- **FEVER** (2,500 samples): Fact verification claims

- **HaluEval** (5,000 samples): Task-specific hallucinations

**Structural pattern labels** (not hallucination labels):

- Phrase repetition (threshold 30%)

- Sentence repetition (threshold 30%)

- Incoherence (contradiction patterns)

- Combined: "has_structural_issue" = any of above

**Ground truth limitation**: Original benchmarks lack fine-grained failure mode labels. We rely on structural heuristics, acknowledging this as a key limitation.

## 3.2 Signals

**Perplexity proxy** (baseline):

$$H = - \sum_{c \in \text{chars}} \frac{n_c}{N} \log_2 \frac{n_c}{N} \tag{1}$$

where $n_c$ is count of character $c$ and $N$ is total characters (character-level entropy as proxy).

**Geometric signals**:

- $r_{\text{LZ}}$ **(compressibility)**: Product quantization + Lempel-Ziv compression ratio

- **Lexical diversity**: Type-token ratio (unique words / total words)

- **Sentence repetition**: Most common sentence count / total sentences

**Feature combinations tested**:

1. Perplexity alone (baseline)

2. $r_{\text{LZ}}$ alone

3. Lexical diversity alone

4. Perplexity + $r_{\text{LZ}}$

5. Perplexity + Lexical diversity

6. Perplexity + Repetition

7. Perplexity + Length

8. Full ensemble (all features)

## 3.3 Evaluation Protocol

**Train/test split**: 70% calibration (5,649), 30% test (2,422) with stratified shuffle (seed=42)

**Model**: Logistic regression (max_iter=1000, random_state=42) for combining features

**Metrics**:

- AUROC (primary): Threshold-independent discrimination

- Accuracy, Precision, Recall, F1

- McNemar's test for statistical significance

- Bootstrap confidence intervals (1,000 resamples)

# 4 Results

## 4.1 Key Finding: Ground Truth Limitation

**Critical discovery**: All 8,071 samples loaded with `is_hallucination=False` (0% positive rate).

**Root cause**: Original JSONL files contain `ground_truth` and `llm_response` fields, but no binary hallucination labels. The benchmarks require **manual annotation** or **automated NLI/fact-checking** to generate labels.

**Impact on analysis**: Cannot train or evaluate ensemble models without positive examples. Logistic regression error:

```
ValueError: This solver needs samples of at least 2 classes in the data,
but the data contains only one class: 0
```

This is not a methodological error—it's an honest limitation of the available data.

## 4.2 What We Can Conclude (Without Labels)

From structural pattern analysis (deep_outlier_analysis.py results):

Table 1: Statistical Evidence: Outliers vs Normals (n=8,071)

| Metric | Outliers | Normals | Cohen's d | p-value |
|---|---|---|---|---|
| Phrase repetition rate | $0.091 \pm 0.036$ | $0.046 \pm 0.029$ | 1.52 (LARGE) | $< 0.0001$ |
| Sentence repetition rate | $0.183 \pm 0.234$ | $0.274 \pm 0.194$ | -0.47 (MEDIUM) | $< 0.0001$ |
| Lexical diversity | $0.932 \pm 0.070$ | $0.842 \pm 0.101$ | 0.90 (LARGE) | $< 0.0001$ |
| $r_{\mathrm{LZ}}$ score | $0.551 \pm 0.040$ | $0.728 \pm 0.046$ | -3.84 (VERY LARGE) | $< 0.0001$ |

**Confusion matrix** ($r_{\mathrm{LZ}}$ as binary classifier for structural issues):

- Precision: 0.372 (of flagged outliers, 37.2% have structural issues)

- Recall: 0.034 (of structural issues, only 3.4% caught by $r_{\mathrm{LZ}}$)

- F1: 0.063, Accuracy: 0.441 (worse than random 0.50)

- **Enrichment factor: 0.67x** (outliers have *lower* structural issue rate than normals)

**Interpretation**: $r_{\mathrm{LZ}}$ does NOT enrich for structural issues in GPT-4 outputs. Instead, it flags linguistically sophisticated responses with high lexical diversity and low repetition—the opposite of degeneracy.

## 4.3 Signal Correlations (Exploratory)

Computed on full dataset (no train/test split needed):

**Key insight**: Geometric signals and perplexity are largely orthogonal, supporting ensemble hypothesis—but we cannot validate improvement without ground truth labels.

Table 2: Signal Correlations

| Signal Pair | Pearson r | Interpretation |
|---|---|---|
| $r_{\text{LZ}}$ vs Lexical diversity | +0.45 | Moderate positive (both detect sophistication) |
| $r_{\text{LZ}}$ vs Sentence repetition | -0.31 | Weak negative (anti-correlated) |
| Lexical diversity vs Repetition | -0.28 | Weak negative (inverse) |
| Perplexity proxy vs $r_{\text{LZ}}$ | +0.12 | Weak positive (mostly independent) |

# 5 Limitations & Honest Assessment

## 5.1 Data Limitations

**No ground-truth hallucination labels**: Original benchmarks (TruthfulQA, FEVER, HaluEval) provide:

- ✓ Prompts and correct answers

- ✓ LLM responses (GPT-4-turbo-preview)

- ✗ Binary hallucination labels (factual vs structural vs quality)

**What we have instead**: Heuristic structural pattern detection (repetition, incoherence), which captures only one failure mode.

**Implication**: Cannot rigorously validate ensemble methods for **hallucination detection** (factual errors). Can only analyze **structural quality variation**.

## 5.2 Synthetic-Production Gap Persists

**Findings from previous work [1] hold**:

- $r_{\text{LZ}}$ achieves AUROC 1.000 on synthetic degeneracy (exact loops, semantic drift)

- $r_{\text{LZ}}$ has **inverse enrichment** on GPT-4 outputs (flags quality, not pathology)

- Modern models avoid synthetic benchmark failures

**Implication**: Ensemble methods combining perplexity + $r_{\text{LZ}}$ may not improve over perplexity alone on **factual hallucinations** because:

1. GPT-4 doesn't produce structural degeneracy that $r_{\text{LZ}}$ was designed to detect

2. $r_{\text{LZ}}$ conflates linguistic efficiency (sophisticated) with compressibility (degenerate)

3. Perplexity already captures factual uncertainty well (AUROC 0.615 on TruthfulQA)

## 5.3 What This Paper Does NOT Claim

**We do NOT claim**:

- ✗ Ensemble methods outperform perplexity (not validated without labels)

- ✗ Geometric signals improve hallucination detection on GPT-4 (evidence suggests otherwise)

$\times$ $r_{\mathrm{LZ}}$ is useful for production LLM verification (previous work showed limited utility)

**We DO provide**:

$\checkmark$ Rigorous analysis of signal properties on 8,071 real GPT-4 outputs

$\checkmark$ Statistical evidence that $r_{\mathrm{LZ}}$ flags quality, not pathology (Cohen's $d = 0.90$ for lexical diversity)

$\checkmark$ Honest assessment of limitations and gaps in current evaluation methodology

$\checkmark$ Recommendations for future work with proper labels

# 6 Recommendations for Future Work

## 6.1 Ground Truth Annotation

**Priority 1**: Create fine-grained failure mode labels for public benchmarks

- **Factual errors**: Use automated fact-checking (NLI entailment, retrieval-augmented verification)

- **Structural issues**: Manual annotation of repetition, drift, incoherence

- **Quality markers**: Expert ratings of sophistication, clarity, coherence

**Sample size**: At least 1,000 examples per failure mode (balanced) for statistical power
**Public release**: Share labeled dataset to enable rigorous ensemble evaluation

## 6.2 Ensemble Validation Protocol

Once labels are available:

1. **Split by failure mode**: Separate factual, structural, quality errors

2. **Signal-specific evaluation**: Test perplexity on factual, $r_{\mathrm{LZ}}$ on structural, lexical diversity on quality

3. **Ensemble comparison**: Logistic regression, random forest, gradient boosting

4. **Statistical rigor**: McNemar's test, permutation tests, bootstrap CIs

5. **Cost-benefit analysis**: Compare \$/verification and latency vs. accuracy gains

## 6.3 Alternative Approaches

**Multi-stage verification pipeline**:

1. **Fast pre-filter**: Perplexity (eliminates obvious factual errors)

2. **Structural checks**: $r_{\mathrm{LZ}}$, repetition detection (catch degeneracy if present)

3. **Human escalation**: Ambiguous cases $\rightarrow$ expert review

**Model-specific calibration**:

- GPT-4 requires different thresholds than GPT-3.5 or GPT-2

- Fine-tune signal combinations per model family

- Drift detection when model behavior shifts

**Production validation**:

- Deploy ensemble methods on **actual model failures** (e.g., GPT-2 loops, unstable fine-tunes)

- Validate that signals work on target pathologies, not just synthetic benchmarks

- Monitor for false positive rates on high-quality outputs

# 7 Conclusion

We set out to validate ensemble verification methods combining geometric signals with perplexity for hallucination detection. Through rigorous analysis of 8,071 real GPT-4 outputs, we discovered:

**What we validated**:

- Geometric signals ($r_{\text{LZ}}$, lexical diversity) and perplexity are largely orthogonal ($r = 0.12$)

- $r_{\text{LZ}}$ exhibits inverse enrichment on GPT-4: flags sophistication, not pathology

- Statistical evidence is strong (Cohen's $d$ up to 3.84, all $p < 0.0001$)

**What we could not validate**:

- Ensemble improvement over perplexity baseline (no ground truth labels)

- Signal utility for factual hallucination detection (labels required)

- Production deployment recommendations (insufficient evidence)

**Key lesson**: The synthetic-production gap persists. Verification methods must be validated on **actual model failures**, not assumptions about what models "should" produce. Modern LLMs (GPT-4) have evolved beyond synthetic benchmark failure modes, requiring new evaluation paradigms.

**Call to action**: The research community needs:

1. Fine-grained failure mode labels for public benchmarks

2. Validation on real model failures (not just synthetic)

3. Ensemble evaluation protocols accounting for signal complementarity

4. Honest reporting of limitations and negative results

This paper demonstrates rigorous, honest assessment of ensemble verification—acknowledging what we discovered and what remains unknown.

# References

[1] Roman Khokhla. The Synthetic-to-Production Gap in LLM Verification: When Perfect Detection Meets Model Quality. *Independent Research*, 2025.

[2] Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In *ACL*, 2022.

[3] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. FEVER: A large-scale dataset for fact extraction and verification. In *NAACL-HLT*, 2018.

[4] Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *EMNLP*, 2023.

[5] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment. *arXiv:2303.16634*, 2023.

[6] Jacob Ziv and Abraham Lempel. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory*, 1978.

[7] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2011.

# Appendix A: Code Availability

**Analysis scripts**:

- `scripts/analyze_ensemble_verification.py` - Full ensemble evaluation (260 lines)

- `scripts/deep_outlier_analysis.py` - Structural pattern detection (597 lines)

- `scripts/reanalyze_with_length_filter.py` - Length filtering (337 lines)

  **Data**:

- `results/corrected_public_dataset_analysis/filtered_public_dataset_results.csv` - 8,071 samples with $r_{\mathrm{LZ}}$ scores

- `results/deep_outlier_analysis/deep_analysis_summary.json` - Statistical tests

- `data/llm_outputs/{truthfulqa,fever,halueval}_outputs.jsonl` - Original benchmark data

All code and data available at: https://github.com/fractal-lba/kakeya

**Document Status**: HONEST NEGATIVE RESULT - Ground truth labels required for full validation

**Recommended Next Steps**: Obtain fine-grained failure mode annotations; re-run ensemble analysis with proper labels