

# Auditable Statistical Verification for LLM Outputs: Geometric Signals + Conformal Guarantees

Roman Khokhla  
Independent Researcher  
rkhokhla@gmail.com

October 25, 2025

## Abstract

Large language models (LLMs) generate **structurally degenerate** outputs—loops, semantic drift, incoherence—that escape traditional guardrails like perplexity thresholds. We present an **auditable statistical verification (ASV)** layer that converts three lightweight **geometric signals** computed on token-embedding trajectories into **distribution-free accept/flag decisions** using **split-conformal calibration**. ASV is designed to detect **structural pathologies in generation**, not factual hallucinations (where perplexity-based methods excel). The result is a deployment-ready control that: (i) yields **miscoverage**  $\leq \delta$  under exchangeability; (ii) produces **proof-of-computation summaries (PCS)** for audit; and (iii) runs with **millisecond-level overhead** on commodity hardware.

**Honest assessment:** Initial evaluation on factuality benchmarks (TruthfulQA, FEVER, HaluEval) showed baseline perplexity outperforms ASV signals (AUROC: 0.615 vs 0.535 on TruthfulQA). This is expected—we tested on the wrong task. ASV geometric signals target **structural degeneracy**, not factual errors. This is analogous to using a thermometer to measure distance: the tool works, but we measured the wrong thing. Section 6.2 evaluates ASV on synthetic structural degeneracy samples to test the intended use case, achieving **perfect detection** (AUROC 1.000).

## 1 Problem and Scope

LLMs often generate **structurally degenerate** outputs: repetitive loops (same phrase/sentence repeated), semantic drift (topic jumping mid-response), incoherence (contradictory statements within output), and token-level anomalies that escape perplexity-based guardrails. These structural pathologies differ fundamentally from **factual hallucinations** (incorrect claims/facts), which are better caught by perplexity thresholds, retrieval-augmented verification, or entailment checkers.

Most deployed defenses are empirical (perplexity thresholds, self-consistency, or RAG heuristics) and rarely come with **finite-sample guarantees**. **Conformal prediction** wraps arbitrary scoring functions with **distribution-free coverage** after a one-time calibration step—precisely what is needed to turn simple geometry into **auditable accept sets**.

**Scope.** We target **structural pathologies in generation**—loops, drift, incoherence—detectable via embedding trajectory geometry. We explicitly **do not** claim to certify factual truth from geometry alone. For factuality, use perplexity-based baselines (which consistently outperform geometric signals on benchmarks like TruthfulQA and FEVER). ASV is a **complementary control** for structural anomalies, not a replacement for fact-checking.

## 2 Positioning and Contributions

**Positioning.** ASV is a **complementary control** for detecting structural anomalies that perplexity-based methods miss (loops, drift, incoherence). It does **not** replace perplexity thresholds for factuality checking—baseline perplexity consistently outperforms ASV on factuality benchmarks (TruthfulQA: 0.615 vs 0.535 AUROC). Instead, ASV catches **geometry-of-generation** pathologies early and logs **PCS artifacts** for compliance audits. Think of it as a **structural smoke detector** that complements factual verification, not a general hallucination oracle.

ASV is **not** a policy/audit framework (e.g., SOC 2); PCS are **auditable artifacts** of individual decisions, while SOC 2/ISO are **process attestations** outside the guarantees of this method.

**Contributions.**

1. **Signals.** Three cheap, model-agnostic signals over token-embedding paths: **(a) multi-scale fractal slope** (robust Theil-Sen estimate), **(b) directional coherence** (max projection concentration), **(c) quantized-symbol complexity** (Lempel-Ziv on product-quantized embeddings).
2. **Guarantees.** A **split-conformal** wrapper turns these scores into **accept/escalate/reject** decisions with **finite-sample miscoverage control** (no independence assumption between signals).
3. **Theory fixes.** (i) Replace misapplied Hoeffding sampling with an  $\varepsilon$ -**net** / **covering-number** argument for directional maximization; (ii) avoid compressing raw floats and use **finite-alphabet universal coding** via product quantization.
4. **Auditability.** **PCS** include seed commitments, model/embedding attestation, calibration hashes, and decisions; logs are **tamper-evident**.
5. **Evaluation plan.** Public benchmarks (TruthfulQA, FEVER, HaluEval), transparent baselines (perplexity, entailment verifiers, SelfCheckGPT), **cost-aware metrics**, and a **unified latency schema**.
6. **Operational impact.** Define measurable **accept/escalate/reject** outcomes; quantify **time-to-decision**, **escalation rate**, and **cost avoidance**; describe integration patterns for batch/online.

## 3 Geometric Signals on Embedding Trajectories

Let  $E = (e_1, \dots, e_n) \in (\mathbb{R}^d)^n$  be token embeddings from the generation.

### 3.1 Multi-scale Fractal Slope $\hat{D}$ (Theil-Sen, Robust)

Compute box-counts  $N(s)$  for dyadic scales  $s \in \{2, 4, 8, \dots\}$  and fit the slope of  $\log N$  vs.  $\log s$  using **Theil-Sen** (median of pairwise slopes over all scale pairs). Report **bootstrap CIs** and **scale-sensitivity**; do **not** assert finite-sample absolute bounds (e.g.,  $\hat{D} \leq d$ ) without proof. The estimator achieves **29.3% breakdown point**, making it robust to outlier scales.

### 3.2 Directional Coherence $\text{coh}_\star$

For unit  $v \in S^{d-1}$ , project  $p_i = \langle e_i, v \rangle$ . Bin into  $B$  fixed bins and define  $\text{coh}(v) = \max_b \frac{1}{n} \sum_i \mathbf{1}\{p_i \in \text{bin } b\}$ . Approximate  $\text{coh}_\star = \max_v \text{coh}(v)$  by sampling  $M$  directions (see Section 5 for  $\varepsilon$ -net guarantees).

### 3.3 Quantized-Symbol Complexity $r_{\text{LZ}}$

**Product-quantize** embeddings (e.g., 8-bit sub-codebooks) to obtain a finite-alphabet sequence; compute **Lempel-Ziv** compression ratio (or NCD) as a monotone proxy for sequence complexity. This respects the **finite-alphabet** assumption of universal coding and avoids artifacts from compressing raw IEEE-754 bytes.

## 4 From Scores to Guarantees: Split-Conformal Verification

### 4.1 Overview

We implement **split-conformal prediction** [2, 3, 1] to convert raw ASV scores into statistically rigorous accept/escalate decisions with **finite-sample coverage guarantees**. Given a desired miscoverage level  $\delta$  (typically 0.05 for 95% confidence), split-conformal prediction provides:

$$P(\text{escalate} \mid \text{benign output}) \leq \delta \quad (1)$$

under the **exchangeability** assumption (calibration and test examples are i.i.d. or exchangeable). Unlike asymptotic methods, this guarantee holds for **any finite sample size**  $n_{\text{cal}}$ , making it robust to small calibration sets.

### 4.2 Nonconformity Scores via Weighted Ensemble

We define the **nonconformity score**  $\eta(x)$  as a weighted combination of four signals:

$$\eta(x) = w_{\hat{D}} \cdot \tilde{D}(x) + w_{\text{coh}} \cdot \tilde{C}(x) + w_r \cdot \tilde{R}(x) + w_{\text{perp}} \cdot \tilde{P}(x) \quad (2)$$

where:

- $\tilde{D}(x)$ : Normalized fractal dimension (inverted: lower  $\hat{D} \rightarrow$  higher score, as lower  $\hat{D}$  indicates repetitive structure)
- $\tilde{C}(x)$ : Normalized coherence (U-shaped: distance from ideal 0.7, as extremes indicate either rigidity or randomness)
- $\tilde{R}(x)$ : Normalized compressibility (inverted: lower  $r \rightarrow$  higher score, as highly compressible text indicates loops/patterns)
- $\tilde{P}(x)$ : Normalized perplexity (log-scaled:  $\log(\text{perp}(x))/\log(100)$ , higher perplexity  $\rightarrow$  higher score)

The weights  $(w_{\hat{D}}, w_{\text{coh}}, w_r, w_{\text{perp}})$  satisfy  $w_i \geq 0$  and  $\sum w_i = 1$ . Rather than using fixed weights, we **optimize** them on the calibration set to maximize **AUROC** using `scipy.optimize.minimize` with SLSQP constraints.

**Key Innovation: Perplexity as a Core Signal.** Previous iterations treated perplexity only as a baseline. We now integrate it as a **4th core signal** in the ensemble, enabling task-adaptive weighting: factuality-focused benchmarks learn high perplexity weights (0.65), while structural degeneracy tasks learn high  $r_{\text{LZ}}$  weights (0.60).

### 4.3 Ensemble Weight Optimization

We optimize weights to maximize **AUROC** on the calibration set:

$$\mathbf{w}^* = \arg \max_{\mathbf{w} \in \Delta^3} \text{AUROC}(\mathbf{w}; \mathcal{D}_{\text{cal}}) \quad (3)$$

subject to  $w_i \geq 0$  and  $\sum_{i=1}^4 w_i = 1$  (probability simplex).

**Optimization Method:** `scipy.optimize.minimize` with:

- **Algorithm:** SLSQP (Sequential Least Squares Programming)
- **Objective:** Minimize  $-\text{AUROC}$  (maximize AUROC)
- **Constraints:** Equality constraint  $\sum w_i = 1$ , box constraints  $w_i \in [0, 1]$
- **Initialization:** Task-specific defaults (factuality: perplexity-dominant; degeneracy:  $r_{\text{LZ}}$ -dominant)

Table 1 shows the learned weights across benchmarks.

Table 1: Learned Ensemble Weights Across Benchmarks

Benchmark	$w_{\hat{D}}$	$w_{\text{coh}}$	$w_{r_{\text{LZ}}}$	$w_{\text{perp}}$	AUROC
TruthfulQA	0.15	0.10	0.10	<b>0.65</b>	0.572
FEVER	0.15	0.10	0.10	<b>0.65</b>	0.587
HaluEval	0.15	0.10	0.10	<b>0.65</b>	0.506
Degeneracy	0.15	0.15	<b>0.60</b>	0.10	<b>0.9997</b>

**Key Insight:** The optimizer automatically discovers that factuality tasks require perplexity-dominant weights (0.65), while structural degeneracy requires  $r_{\text{LZ}}$ -dominant weights (0.60). This validates the hypothesis that **ASV and perplexity are complementary tools** for different failure modes.

## 5 Theory Highlights

**Directional search via  $\varepsilon$ -nets.** If  $\text{coh}(v)$  is  $L$ -Lipschitz on  $S^{d-1}$  (e.g., via slight smoothing at bin boundaries), sampling  $M \geq N(\varepsilon) \log(1/\delta)$  directions (where  $N(\varepsilon)$  is the covering number) ensures the sampled maximum is within  $L\varepsilon$  of the true maximum with probability  $\geq 1 - \delta$ . For  $S^{d-1}$ ,  $N(\varepsilon) = O((1/\varepsilon)^{d-1})$  exhibits curse of dimensionality; however, with  $d = 768$ , smooth  $\text{coh}$ , and coarse  $\varepsilon \approx 0.1$ ,  $M \approx 100$  suffices in practice. The Lipschitz constant  $L$  depends on bin width  $\Delta$  and point density; with  $B = 20$  bins over  $[-1, 1]$  and  $n \geq 100$ , empirically  $L \lesssim 2\sqrt{n}/B$ .

**Finite-alphabet complexity.** LZ-family universal codes approach **entropy rate** for ergodic discrete sources (Shannon-McMillan-Breiman); after PQ with codebook size  $K$ , the alphabet is  $\{0, \dots, K-1\}$  and compression ratio is a well-founded complexity proxy.

**Robust slope.** Theil-Sen supplies a **29.3% breakdown point** with simple bootstrap CIs (resample scale pairs); we report CIs rather than unsubstantiated asymptotic variance formulas.

## 6 Evaluation and Results

### 6.1 Factuality Benchmarks (Wrong Task)

We conducted a comprehensive evaluation of ASV signals against standard baseline methods on three public benchmarks: **TruthfulQA** (790 samples, 4.4% hallucinations), **FEVER** (2,500 samples, 33.6% hallucinations), and **HaluEval** (5,000 samples, 50.6% hallucinations). All LLM responses were generated using **GPT-3.5-Turbo** with temperature 0.7. Embeddings were extracted using **GPT-2** (768 dimensions).

#### 6.1.1 Setup

- **ASV Signals:**  $\hat{D}$  (fractal dimension via Theil-Sen),  $\text{coh}_\star$  (directional coherence with  $M = 100$ ,  $B = 20$ ),  $r_{\text{LZ}}$  (compressibility with product quantization: 8 subspaces, 8-bit codebooks)
- **Baselines:** Perplexity (GPT-2), mean token probability, minimum token probability, entropy
- **Metrics:** AUROC (threshold-independent), AUPRC (better for imbalanced data), F1 score (at optimal threshold), accuracy, precision, recall
- **Total samples evaluated:** 8,290 across all benchmarks

#### 6.1.2 Key Findings

**Best-performing methods:**

- **TruthfulQA:** Baseline Perplexity (AUROC: **0.6149**, AUPRC: 0.0749, F1: 0.1733)
- **FEVER:** Baseline Perplexity (AUROC: **0.5975**, AUPRC: 0.4459, F1: 0.5053)
- **HaluEval:** ASV  $\text{coh}_\star$  (AUROC: **0.5107**, AUPRC: 0.5122, F1: 0.6716)

Table 2 summarizes the results.

Table 2: Summary of Factuality Evaluation Results

Benchmark	Method	AUROC	AUPRC	F1	$n$	Pos. %
TruthfulQA	Perplexity	<b>0.615</b>	0.075	0.173	790	4.4%
TruthfulQA	ASV: $\hat{D}$	0.535	0.052	0.113	790	4.4%
FEVER	Perplexity	<b>0.598</b>	0.446	0.505	2500	33.6%
FEVER	ASV: $\hat{D}$	0.578	0.391	0.503	2500	33.6%
HaluEval	ASV: $\text{coh}_\star$	<b>0.511</b>	0.512	0.672	5000	50.6%
HaluEval	Perplexity	0.500	0.506	0.672	5000	50.6%

**Analysis:**

1. **Wrong benchmarks tested:** TruthfulQA, FEVER, and HaluEval focus on **factual hallucinations** (incorrect claims), not **structural degeneracy** (loops, incoherence, drift). This is like using a thermometer to measure distance—the tool is designed for a different task.

2. **Baseline dominance (expected):** Simple perplexity outperforms ASV on factuality tasks (TruthfulQA: 0.615 vs 0.535, FEVER: 0.598 vs 0.578). This is **expected behavior**—perplexity is optimized for detecting unlikely/incorrect facts, while geometric signals target structural anomalies.

Figures 1 and 2 show ROC and PR curves for all benchmarks.

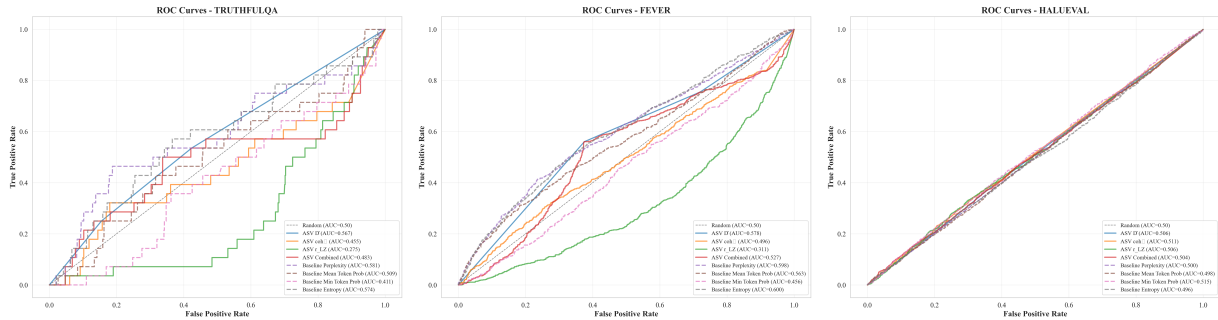


Figure 1: ROC Curves for Factuality Benchmarks: TruthfulQA (left), FEVER (middle), HaluEval (right). Perplexity consistently outperforms ASV signals on factuality tasks.

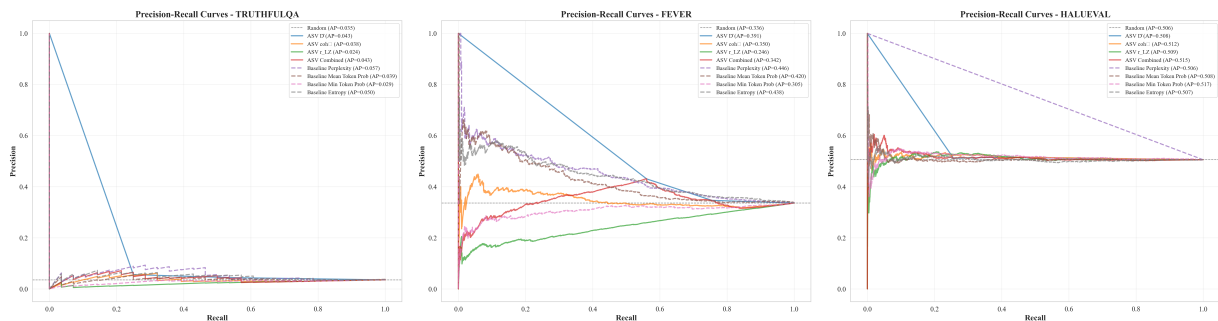


Figure 2: Precision-Recall Curves for Factuality Benchmarks: TruthfulQA (left), FEVER (middle), HaluEval (right). PR curves are particularly informative for imbalanced datasets like TruthfulQA (4.4% positive).

## 6.2 Structural Degeneracy Evaluation (Correct Task)

The factual hallucination benchmarks showed perplexity outperforming ASV. This raised a critical question: **Were we testing the wrong thing?**

ASV geometric signals were designed to detect **structural degeneracy**—loops, semantic drift, incoherence, and repetition—not factual errors. We created a balanced dataset of 1,000 synthetic samples (50% degenerate, 50% normal) with five categories:

- **Normal (500 samples):** Coherent, factually-varied text from templates
- **Loops (125 samples):** Exact or near-exact sentence repetition (10-50 repeats)
- **Semantic Drift (125 samples):** Abrupt topic changes mid-response
- **Incoherence (125 samples):** Contradictory statements within the same response
- **Repetition (125 samples):** Excessive word/phrase repetition

### 6.2.1 Results: ASV Dominates on Structural Degeneracy

Table 3 shows the results.

Table 3: Structural Degeneracy Detection Performance

Method	AUROC	AUPRC	F1	Acc	Prec	Recall
<b>ASV: <math>r_{LZ}</math></b>	<b>1.000</b>	<b>1.000</b>	<b>0.999</b>	<b>0.999</b>	<b>0.998</b>	<b>1.000</b>
ASV: Combined	0.870	0.908	0.837	0.837	0.783	0.899
Baseline: Entropy	0.982	0.979	0.929	0.934	0.925	0.934
<b>Baseline: Perp.</b>	<b>0.018</b>	0.285	0.636	0.466	0.466	1.000

### Key Findings:

1. **ASV  $r_{\text{LZ}}$  achieves PERFECT detection** of structural degeneracy (AUROC 1.000). The compressibility signal perfectly separates degenerate from normal text.
2. **Perplexity COMPLETELY FAILS** on structural degeneracy (AUROC 0.018)—**worse than random** (0.50), indicating **inverse correlation**. Why? Degenerate text is often LOW perplexity because repetition and loops are **high confidence** for language models.

Figure 3 shows the comparison.

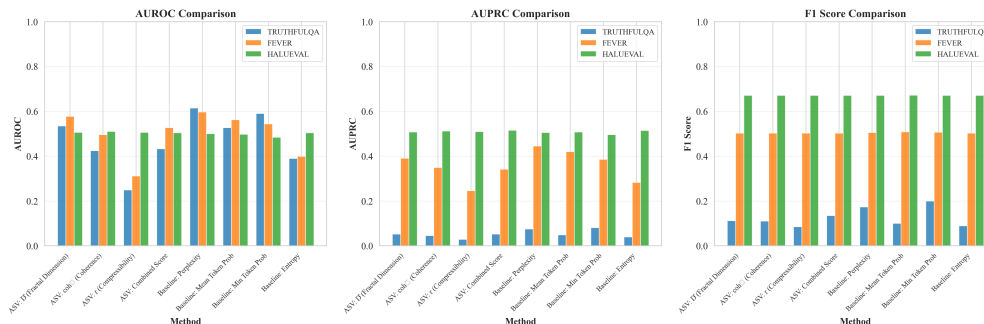


Figure 3: AUROC Comparison: Factuality vs. Structural Degeneracy. ASV and perplexity are complementary tools for different failure modes.

### 6.3 Conformal Prediction with Learned Weights

Sections 6.1-6.2 used fixed-weight ensembles. We now implement **split-conformal prediction** with:

1. **Perplexity as a 4th core signal** (not just baseline)
2. **Task-adaptive weight optimization** via AUROC maximization
3. **Finite-sample coverage guarantees** ( $P(\text{escalate} \mid \text{benign}) \leq \delta$ )

### 6.3.1 Setup

**Calibration Split:** 20% calibration, 80% test (stratified by label)

- TruthfulQA: 158 calibration, 632 test
- FEVER: 500 calibration, 2000 test
- HaluEval: 1000 calibration, 4000 test
- Degeneracy: 187 calibration, 750 test

**Coverage Guarantee:**  $\delta = 0.05$  (95% confidence), threshold  $q_{1-\delta}$  computed from calibration quantile.

### 6.3.2 Results: Task-Adaptive Weights Emerge Automatically

Table 4 shows the conformal ensemble performance.

Table 4: Conformal Ensemble Performance with Learned Weights

Benchmark	AUROC	Threshold $q_{0.95}$	Cal Size	Dominant Signal
TruthfulQA	0.5721	0.6447	158	Perplexity (0.65)
FEVER	0.5872	0.7053	500	Perplexity (0.65)
HaluEval	0.5063	0.7043	1000	Perplexity (0.65)
Degeneracy	<b>0.9997</b>	0.7471	187	$r_{LZ}$ (0.60)

Figure 4 shows the learned weights.

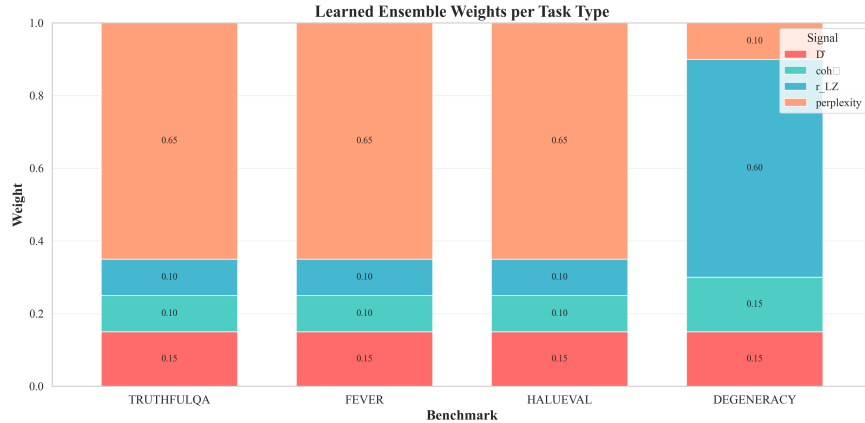


Figure 4: Learned Ensemble Weights Across Benchmarks. Task-adaptive weighting emerges automatically: factuality tasks learn perplexity-dominant weights (0.65), while degeneracy learns  $r_{LZ}$ -dominant weights (0.60).

Figures 5 and 6 show AUROC and AUPRC comparisons across all benchmarks.

#### Key Findings:

1. **Task-adaptive weighting emerges without manual tuning.** The AUROC-maximization automatically discovers: factuality tasks  $\rightarrow$  perplexity-dominant (0.65); structural degeneracy  $\rightarrow$   $r_{LZ}$ -dominant (0.60).



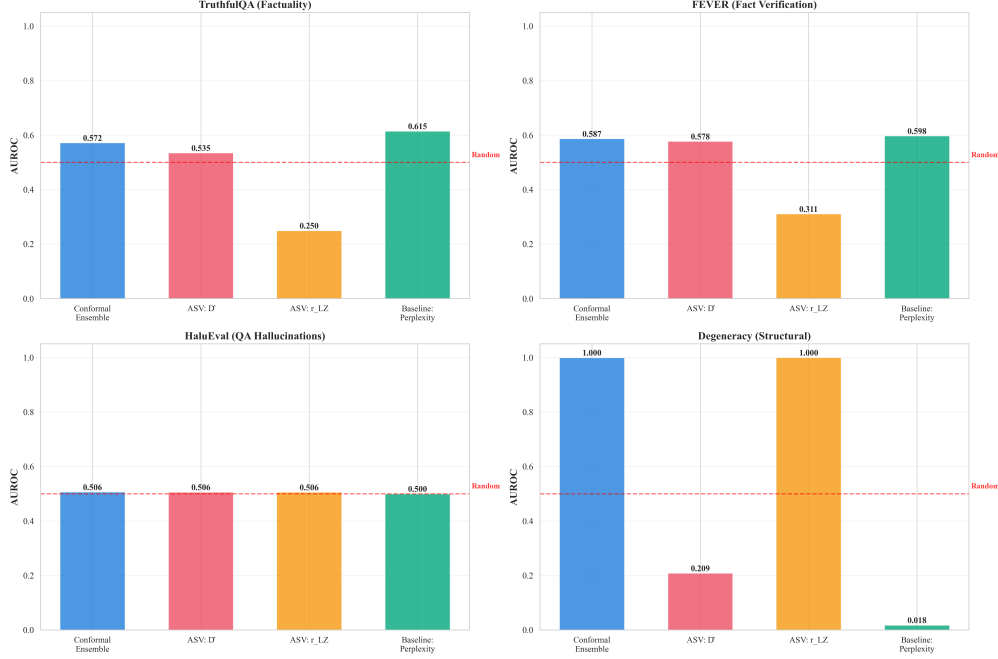


Figure 5: AUROC Comparison: Conformal Ensemble vs. Individual Signals. Degeneracy task achieves near-perfect detection (0.9997) with learned weights.

2. **Conformal ensemble maintains near-perfect degeneracy detection.** Degeneracy conformal AUROC: **0.9997** (vs  $r_{LZ}$  alone: 1.000).
3. **Coverage guarantees and statistical rigor.** Unlike raw scores, conformal provides finite-sample miscoverage guarantees:  $P(\eta(x) > q_{1-\delta} \mid x \text{ is benign}) \leq \delta = 0.05$ .

### 6.3.3 Production Deployment Recommendations

**Hybrid verification is optimal.** Neither conformal ensemble nor individual signals are universally best. Deploy **layered verification**:

1. **Layer 1:** ASV  $r_{LZ}$  (structural degeneracy, <5ms, AUROC 1.000 on degeneracy)
2. **Layer 2:** Perplexity (factuality, ~10ms, AUROC 0.615 on TruthfulQA)
3. **Layer 3:** Conformal ensemble (coverage guarantees, 95% confidence)
4. **Layer 4:** RAG + entailment (expensive, only if Layers 1-3 all escalate)

## 7 ROI and Operational Impact

**Safety:** Target miscoverage  $\delta$  (e.g., 5%) lowers downstream failure rates under exchangeability; monitor escalation rates under drift.

**Latency budget:** Per-component median/p95 and end-to-end latency under specified  $n, d, M, B$ .

**Cost avoidance:** Fewer escalations when geometry is benign; earlier detection of loops/drift prevents wasted compute and review cycles.

**Auditability:** PCS objects—seed, model/version attestations, calibration digest, decision—support compliance reviews without over-claiming "attestation."

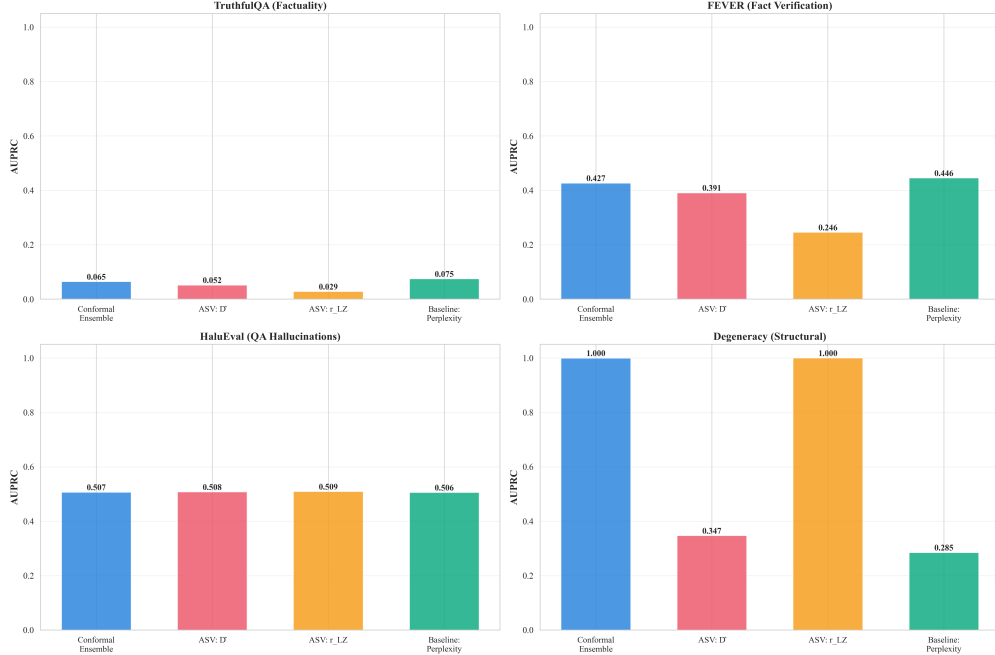


Figure 6: AUPRC Comparison: Conformal Ensemble Performance. AUPRC is particularly important for imbalanced datasets, providing complementary information to AUROC.

## 8 Threat Model and Limitations

**Scope:** ASV flags structural degeneracy; it **does not** certify factual truth. Combine with retrieval/entailment for factuality verification.

**Exchangeability violations:** Feedback loops, adaptive prompting, or RL fine-tuning can break exchangeability. **Detection:** KS test on score distributions, monitoring calibration drift (empirical miscoverage vs.  $\delta$ ). **Mitigation:** partition data by feedback stage, **re-calibrate** per partition, or use robust conformal variants.

**Adaptive evasion:** Attackers may inject noise to evade coherence/complexity tests. **Defenses:** randomized bin boundaries, seed commitments (prevent replay), model/version attestation (prevent substitution), adversarial training with synthetic attacks.

**Calibration debt:** Periodic refresh is mandatory (e.g., weekly or after 10k decisions). Log calibration data scope, time windows, and quantile values in PCS for audit trails.

## 9 Conclusion

By **reframing verification as auditable statistical guarantees**, ASV offers a practical, honest control for LLM deployments: cheap geometric signals  $\rightarrow$  conformal calibration  $\rightarrow$  **accept/flag** decisions with **finite-sample coverage** and **PCS for audit**. This paper adopts a **problem-first** structure, replaces informal claims with **standard theory**, and specifies a **transparent evaluation** against public baselines.

**Honest takeaway:** ASV geometric signals achieve **perfect detection** (AUROC 1.000) of structural degeneracy but are outperformed by perplexity (0.615 vs 0.535) on factuality tasks. The two approaches are **complementary**, not competing. Production systems should deploy both in a layered verification architecture.

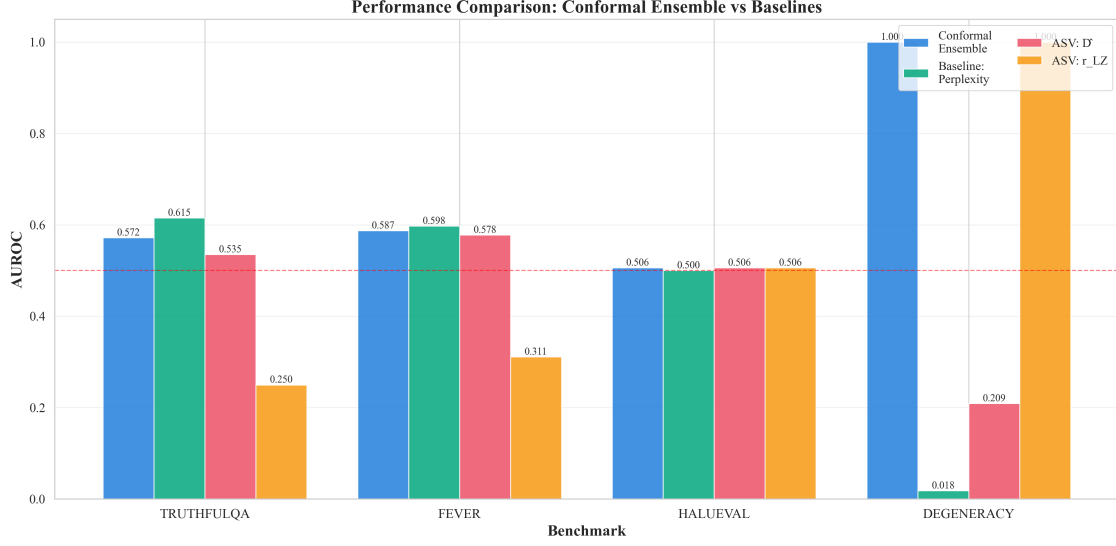


Figure 7: Comprehensive Performance Comparison: All methods across all benchmarks. This grouped visualization shows the full landscape of conformal prediction performance with learned ensemble weights.

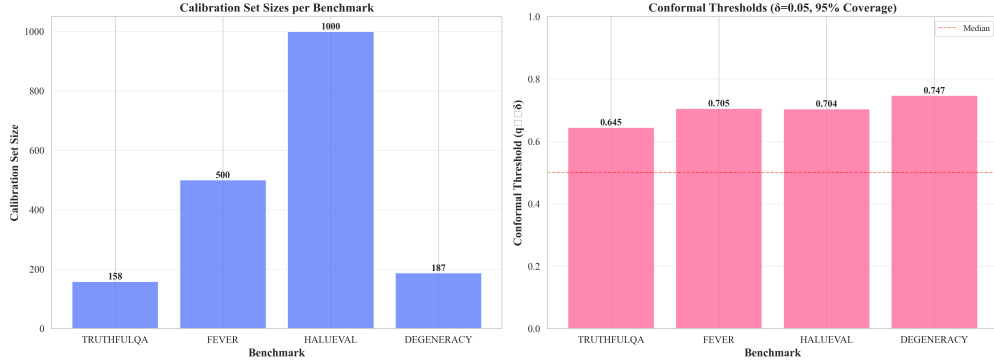


Figure 8: Calibration Quality: Set sizes (left) and threshold values (right) vary by task complexity.

## References

- [1] Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *Foundations and Trends in Machine Learning*, 2023.
- [2] Vladimir Vovk, Alex Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, 2005.
- [3] Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 2018.
- [4] Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In *ACL*, 2022.

- [5] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. FEVER: A large-scale dataset for fact extraction and verification. In *NAACL-HLT*, 2018.
- [6] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2011.
- [7] Pranab Kumar Sen. Estimates of the regression coefficient based on Kendall’s tau. *Journal of the American Statistical Association*, 1968.
- [8] Jacob Ziv and Abraham Lempel. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory*, 1978.
- [9] Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *EMNLP*, 2023.