

TriLens: Per-Layer Logit-Lens Entropy for White-Box Hallucination Detection

Anonymous ACL submission

Abstract

When a language model hallucinates, the final answer is wrong, but the mistake is not necessarily invisible inside the model. Different internal pathways may remain uncertain, disagree in how quickly they sharpen, or commit to competing continuations before the output is produced. We introduce *TriLens*, a white-box detector that turns this intuition into a compact representation: at every layer, it reads the multi-head self-attention output, the feed-forward output, and the residual stream through the model’s own logit lens, then records only the entropy of each readout. The resulting $3L$ -dimensional trajectory describes how certainty forms across depth and across modules, without storing high-dimensional hidden states or sampling multiple generations. This simple signal yields a strong detector across instruction-tuned LLMs and QA benchmarks, and our analyses show that the three module-wise entropy trajectories provide complementary evidence. *TriLens* suggests that hallucination detection can benefit from tracking how internal computation settles, not only what the final layer predicts. The code is available at <https://anonymous.4open.science/w/TriLens/>.

1 Introduction

Large language models (LLMs) have reshaped a wide range of NLP tasks, but their tendency to generate hallucinations—outputs that are fluent yet factually wrong or inconsistent with provided context—remains a fundamental obstacle to deployment in knowledge-sensitive settings (Ji et al., 2023; Li et al., 2023). Reliable *detection* enables selective abstention, confidence calibration, and targeted retrieval.

Among existing detection approaches, *white-box* methods that analyze an LLM’s internal state during a single forward pass are particularly attractive: they avoid the need for multiple samples or external references while enabling efficient

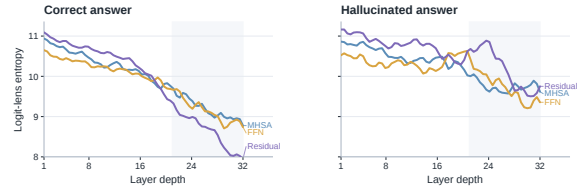


Figure 1: **TriLens: Per-layer logit-lens entropy tracks internal certainty during generation.** Supported answers show coordinated entropy sharpening across MHSA, FFN, and residual-stream readouts, whereas hallucinated answers tend to retain higher, less stable, and less synchronized entropy across depth.

inference-time scoring. Recent work has shown that hallucination-relevant information can be extracted from hidden states, gradients, semantic-entropy probes, spectral summaries, attention kernels, and update-consistency features (Azaria and Mitchell, 2023; Kossen et al., 2024; Hu et al., 2024; Chen et al., 2024; Sriramanan et al., 2024; Zhang et al., 2025). However, these approaches still largely treat internal computation as a high-dimensional representation to be classified, rather than asking which low-dimensional aspect of the computation changes when a model moves toward an unsupported answer.

This leads to two empirical questions. First, can hallucination be detected from how internal certainty forms across depth, rather than from final-layer confidence or high-dimensional hidden states? Second, do the main transformer pathways provide redundant views of that uncertainty, or do attention, feed-forward, and residual-stream readouts carry complementary evidence? A useful detector should answer these questions under a single forward pass, expose where uncertainty enters the computation, and avoid turning every hidden state into a large learned representation. We therefore test a simpler alternative: instead of storing activations themselves, we ask how uncertain different parts of the computation look when read through

072	the model’s own vocabulary lens.	
073	<i>TriLens</i> records one scalar from each of three in-	
074	ternal locations at every layer: the multi-head self-	
075	attention output, the feed-forward output, and the	
076	residual stream. Each location is projected through	
077	the logit lens, converted to a Shannon entropy, and	
078	concatenated across depth. The result is a compact	
079	3 <i>L</i> -dimensional trajectory that tracks how certainty	
080	forms across pathways. It is small enough for a	
081	simple probe, but structured enough to distinguish	
082	whether contextual routing, parametric recall, and	
083	the composed residual state sharpen together or	
084	remain diffuse.	
085	Why should a first-order statistic suffice? Our	
086	working hypothesis is that hallucination-relevant	
087	errors are often accompanied by elevated uncer-	
088	tainty across internal pathways of an LLM, and	
089	that per-layer Shannon entropy of logit-lens distri-	
090	butions is a compact diagnostic correlate of this	
091	effect. In a decoder transformer, MHSA routes	
092	contextual information while the FFN pathway re-	
093	trieves parametric content (Geva et al., 2021, 2023;	
094	Elhage et al., 2021). When these pathways align,	
095	their logit-lens distributions typically sharpen with	
096	depth; when they diverge, refinement can become	
097	slower and higher-entropy. <i>TriLens</i> tracks this be-	
098	havior with three scalars per layer, yielding a com-	
099	compact alternative to high-dimensional hidden-state	
100	features.	
101	Contributions. Our main contribution is to test	
102	these questions through a compact pathway-wise	
103	entropy feature, supported by ablations and diag-	
104	nostic comparisons.	
105	• We introduce TriLens , a per-layer logit-lens en-	
106	trophy feature over MHSA, FFN, and residual-	
107	stream states.	
108	• Across three instruction-tuned LLMs and four	
109	benchmark settings, this 3 <i>L</i> -dimensional feature	
110	improves over the strong prior white-box base-	
111	line ICR Probe on <i>all 12</i> cells, averaging +12.1	
112	AUROC, and attaining the highest AUROC in 11	
113	of the 12 cells.	
114	• Ablations show that MHSA, FFN, and residual-	
115	stream entropies contribute complementary sig-	
116	nal, while an intra-layer module-disagreement	
117	term adds little.	
118	• A head-to-head comparison with DoLa-style	
119	cross-layer contrast features shows that per-layer	
120	entropy captures much of the useful signal in that	
	family.	121
	• Layer-wise analysis shows that the peak-	122
	discriminative depth varies systematically across	123
	architectures and benchmarks, consistent with	124
	differences in how models integrate contextual	125
	and parametric information.	126
	2 Related Work	127
	White-box hallucination detection. The dom-	128
	inant white-box line trains detectors on in-	129
	ternal activations from a single forward pass.	130
	SAPLMA (Azaria and Mitchell, 2023), layer-	131
	combination probes (CH-Wang et al., 2024),	132
	SEP (Kossen et al., 2024), gradient-aware detec-	133
	tors (Hu et al., 2024), INSIDE (Chen et al., 2024),	134
	and LLM-Check (Sriramanan et al., 2024) all	135
	fit this template, while ICR Probe (Zhang et al.,	136
	2025) reframes detection around residual-stream	137
	updates. <i>TriLens</i> differs by using a compact 3 <i>L</i> -	138
	scalar feature built from pathway-specific logit-lens	139
	entropy rather than high-dimensional hidden states	140
	or update-consistency scores.	141
	Uncertainty-based detection. A parallel black-	142
	box line scores hallucination risk from a model’s	143
	own uncertainty, including calibration-style sig-	144
	nals (Kadavath et al., 2022; Malinin and Gales,	145
	2020), semantic uncertainty and semantic en-	146
	tropy (Kuhn et al., 2023; Farquhar et al., 2024),	147
	and recent refinements based on Bayesian estima-	148
	tion, pairwise semantic similarity, fact-level self-	149
	consistency, or log-probability time series (Ciosek	150
	et al., 2025; Nguyen et al., 2025; Sawczyn et al.,	151
	2026; Shapiro et al., 2026). <i>TriLens</i> instead uses	152
	a single white-box pass over internal activations,	153
	without multiple sampled generations or output-	154
	layer uncertainty alone.	155
	Logit-lens and mechanistic motivation. The	156
	logit lens and Tuned Lens project intermediate	157
	states into vocabulary space (nostalgebraist, 2020;	158
	Belrose et al., 2023); DoLa (Chuang et al., 2024)	159
	contrasts these distributions across layers. Our de-	160
	sign is further motivated by mechanistic views of	161
	FFNs as key-value memory, MHSA as residual-	162
	stream routing, and depth-varying module influ-	163
	ence (Geva et al., 2021, 2023; Elhage et al., 2021;	164
	Stolfo et al., 2023). <i>TriLens</i> combines these in-	165
	gredients into a supervised detection feature based	166
	on pathway-specific entropy trajectories; §4.6 com-	167
	pares it directly with DoLa-style cross-layer con-	168
	trast.	169

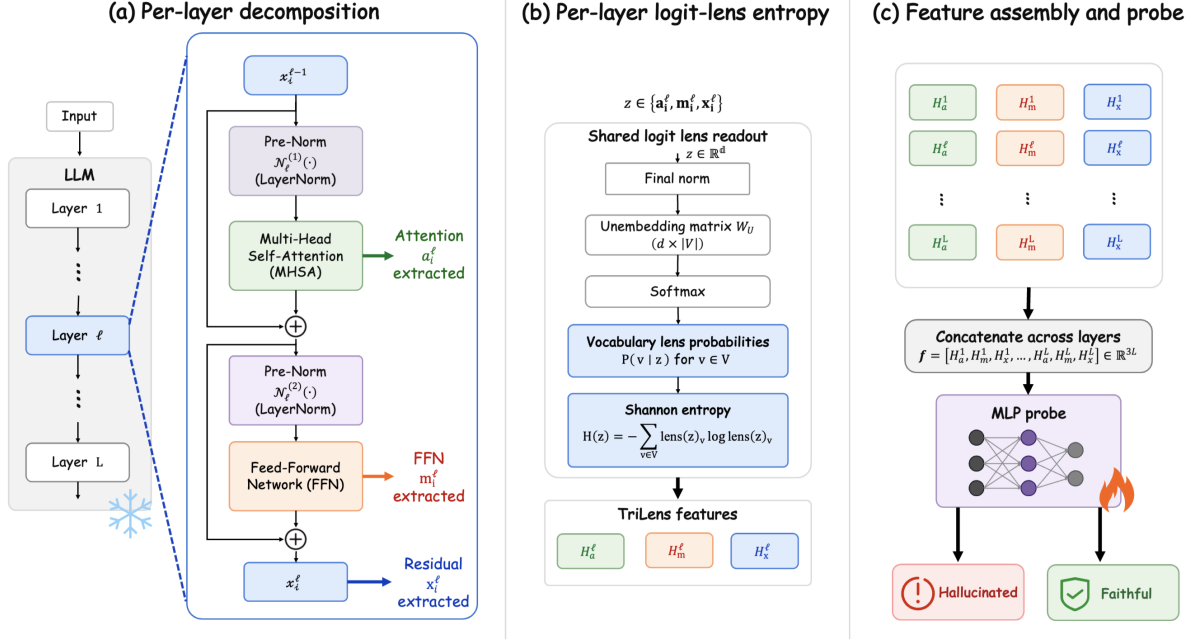


Figure 2: TriLens overview. A single forward pass over the evaluation sequence extracts the MHSA output, FFN output, and residual stream at every layer; each is projected through the model’s logit lens, converted to a Shannon entropy, concatenated into a $3L$ -dimensional feature, and fed to a probe.

3 Method

3.1 Per-Module Logit-Lens Entropy

We consider an LLM with L decoder layers and vocabulary V . At layer ℓ and token position i , the residual stream updates as

$$x_i^\ell = x_i^{\ell-1} + \underbrace{\text{MHSA}^\ell(\mathcal{N}_\ell(x_i^{\ell-1}))}_{a_i^\ell} + \underbrace{\text{FFN}^\ell(\mathcal{N}_\ell(x_i^{\ell-1} + a_i^\ell))}_{m_i^\ell}. \quad (1)$$

For a hidden activation $z \in \{x_i^\ell, a_i^\ell, m_i^\ell\}$, the *logit lens* (nostalgebraist, 2020) projects z into vocabulary-space probabilities

$$\text{lens}(z) = \text{softmax}(W_U \text{Norm}_{\text{final}}(z)), \quad (2)$$

$$H(z) = - \sum_{v \in V} \text{lens}(z)_v \log \text{lens}(z)_v.$$

using the model’s own unembedding matrix W_U and final normalization layer. Here \mathcal{N}_ℓ denotes the layer’s pre-submodule normalization (e.g., RMSNorm in the model families studied here), while the lens itself always uses the model’s *final* normalization layer before unembedding. This mirrors the standard logit-lens construction on the residual stream and keeps a^ℓ , m^ℓ , and x^ℓ in the same read-out coordinate. In implementation we apply the

final normalization directly to each of a^ℓ , m^ℓ , and x^ℓ , with no additional centering or rescaling. The *per-layer logit-lens entropy* is the Shannon entropy of this distribution; we compute it at all three positions, yielding $H_a^\ell, H_m^\ell, H_x^\ell$ for every layer (Eq. 2). TriLens applies the same readout not only to the residual stream but also to the isolated MHSA and FFN writes, using the resulting three-pathway trajectory as a detection feature.

3.2 Mechanistic Motivation

The choice to measure entropy *per module* rather than only on x^ℓ , and to use *Shannon entropy* rather than richer summary statistics, is not arbitrary. It follows from three established views of transformer dynamics.

Dual-pathway decomposition. Under the residual-stream framework of Elhage et al. (2021), a^ℓ and m^ℓ are functionally distinct pathways: the OV circuit of MHSA copies content from earlier positions, while FFN functions as a read-out from a parametric key-value memory (Geva et al., 2021, 2023). In factual completion, both pathways often converge on the same target token—MHSA because the context uniquely determines the answer, FFN because the parametric memory has stored the relevant key–value pair. In hallucinated completion under our benchmark setting, the pathways can

216 *diverge*: a^ℓ may still route contextually correct
 217 information, while m^ℓ can partially recall the same
 218 correct alternative even though the observed token
 219 is wrong. Projecting a^ℓ and m^ℓ independently
 220 through the logit lens exposes this divergence,
 221 which a single measurement on the composed x^ℓ
 222 can mask whenever the two pathways cancel in
 223 direction but not in magnitude. Although a^ℓ and
 224 m^ℓ are not themselves full residual states, they
 225 are additive writes in the same residual space and
 226 are therefore legitimate objects for the model’s
 227 own final readout. The resulting distributions
 228 should be interpreted as *readout preferences*
 229 *of an isolated write* rather than as standalone
 230 next-token predictions; our claim is empirical and
 231 diagnostic, not that a^ℓ or m^ℓ alone constitute a full
 232 decoder state. Accordingly, the role of H_a^ℓ and
 233 H_m^ℓ in TriLens is not to replace H_x^ℓ but to supply
 234 complementary pathway-specific readouts under a
 235 shared lens.

236 **Superposition and commitment.** The residual
 237 stream at any layer is a linear superposition of all
 238 prior module writes. When the writes are coherent
 239 on a single vocabulary direction, the lens distribu-
 240 tion is sharply peaked; when they are incoherent,
 241 the distribution spreads over multiple candidates.
 242 Shannon entropy is a compact statistic that tracks
 243 this spreading: it is invariant to which candidate
 244 tokens are involved, makes no assumption about
 245 whether the vocabulary is flat or structured, and is
 246 a calibrated function of the distribution’s concen-
 247 tration. Richer alternatives such as KL divergence
 248 to the final layer (DoLa’s signal) or a spectral sum-
 249 mary lose information about within-layer commit-
 250 ment, as we show empirically in §4.6.

251 **Iterative refinement trajectories.** The logit-lens
 252 view treats each layer as producing a refinement
 253 of a running next-token prediction (nostalgabraist,
 254 2020; Belrose et al., 2023). Under this view,
 255 $H(\text{lens}(z^\ell))$ is the uncertainty of the model’s be-
 256 lief at depth ℓ along that module’s pathway. For a
 257 correct continuation, uncertainty often decreases
 258 with depth as module contributions reinforce each
 259 other. For a hallucinated continuation, the refine-
 260 ment trajectory is more often non-monotone: un-
 261 certainty can re-enter the distribution at the depth
 262 where the parametric pathway recalls the true to-
 263 ken while the forced context fixes the output on the
 264 wrong one. The per-layer vector $(H_a^\ell, H_m^\ell, H_x^\ell)_\ell$
 265 can thus be viewed as a trajectory of this three-
 266 pathway refinement process, and a linear classifier

267 on this trajectory can separate the two regimes in
 268 our benchmark setting.

269 These three views explain why the feature is
 270 minimal (three location-specific scalars per layer),
 271 why it is non-redundant (each scalar measures a
 272 different pathway’s commitment), and why it is ar-
 273 chitecturally natural (it uses only the model’s own
 274 unembedding and pre-trained norm). Section 4.4
 275 examines the resulting empirical layer-wise trajec-
 276 tories.

277 3.3 Feature Construction

278 Given an evaluation sequence, we run the LLM
 279 once and extract $(H_a^\ell, H_m^\ell, H_x^\ell)$ at the token po-
 280 sitions used for scoring. Following the standard
 281 probe-input convention, we obtain a fixed-length
 282 vector per sample by applying a fixed readout rule
 283 consistently throughout the paper. The resulting
 284 feature is

$$285 \mathbf{f} = (H_a^1, H_m^1, H_x^1, \dots, H_a^L, H_m^L, H_x^L) \in \mathbb{R}^{3L}. \quad (3)$$

286 For the models we consider, $3L$ ranges from
 287 84 (Qwen2.5-7B, $L=28$) to 126 (Gemma-2-9B,
 288 $L=42$)—roughly $30\times$ smaller than the raw hidden-
 289 state features used by SAPLMA and its descen-
 290 dants.

291 3.4 Probes

292 We evaluate two classifiers on \mathbf{f} : (i) a **linear probe**
 293 (L2-regularized logistic regression), which high-
 294 lights that our feature is already separable without
 295 a learned transformation; and (ii) an **MLP probe**
 296 ($3L \rightarrow 128 \rightarrow 64 \rightarrow 32 \rightarrow 1$; LeakyReLU(0.01),
 297 BatchNorm, Dropout 0.3), using the same hidden-
 298 layer topology as the architecture used by the
 299 strong prior baseline (Zhang et al., 2025). Because
 300 TriLens uses a $3L$ -dimensional input whereas ICR
 301 Probe uses an L -dimensional input, this compar-
 302 ison controls probe family and training recipe
 303 but does not fully equalize input dimensionality.
 304 We therefore treat the MLP comparison as
 305 a matched probe-family comparison rather than
 306 as a dimension-controlled capacity study, and we
 307 also report linear-probe results plus feature-subset
 308 ablations to separate feature quality from probe
 309 size. As an additional dimension-matched con-
 310 trol, Appendix B compares TriLens against a
 311 repeated-feature baseline $H_x^{\times 3}$ that simply copies
 312 the residual-stream entropy three times to match
 313 the same $3L$ input size without adding new infor-
 314 mation, and further reports stricter MLP controls

that reduce the three-pathway feature back to L dimensions via simple summaries or train-split PCA. Training details (BCE loss, Adam, 50 epochs, LR scheduler) are identical to those in Zhang et al. (2025) and are listed in Appendix H.

4 Experiments

4.1 Experimental Setup

Models. We evaluate on three recent instruction-tuned LLMs spanning the $\sim 7\text{--}9\text{B}$ parameter range and three architecture families: Qwen2.5-7B-Instruct (Yang et al., 2024), Meta-Llama-3-8B-Instruct (Grattafiori et al., 2024), and Gemma-2-9B-it (Team et al., 2024). The first two use the standard RMSNorm+GQA decoder stack; Gemma-2 alternates full and sliding-window attention.

Datasets. We evaluate on four benchmarks commonly used in prior white-box detection work: HaluEval-QA (Li et al., 2023), SQuAD2.0 (Rajpurkar et al., 2018), HotpotQA (Yang et al., 2018), and TriviaQA (Joshi et al., 2017). For datasets that natively provide positive/negative supervision (HaluEval, SQuAD2), we use the provided supervision directly. For HotpotQA and TriviaQA, we follow a fixed benchmark-specific preprocessing pipeline to instantiate the evaluation examples used in our experiments. The same preprocessed instances are shared by TriLens and all baselines.

Metrics and protocol. Each benchmark is split 80/20 into train and test sets, and we report test-set AUROC averaged over five random seeds. Unless otherwise stated, we randomly sample 10,000 instances from each dataset, keeping the final label distribution balanced within every dataset. All model states are extracted from a single forward pass over the evaluation sequence, and the same sampled instances, splits, model suite, and model-dataset grid are shared by TriLens and all baselines. Further implementation details are listed in Appendix H.

Baselines. We compare against six baselines spanning both training-free and trainable white-box detectors. The training-free baselines are perplexity (PPL) (Ren et al., 2022), length-normalized entropy (LN-Entropy) (Malinin and Gales, 2020), and LLM-Check (Sriramanan et al., 2024); the trainable baselines are SAPLMA (Azaria and Mitchell, 2023), SEP (Kossen et al., 2024), and the previous state-of-the-art ICR Probe (Zhang

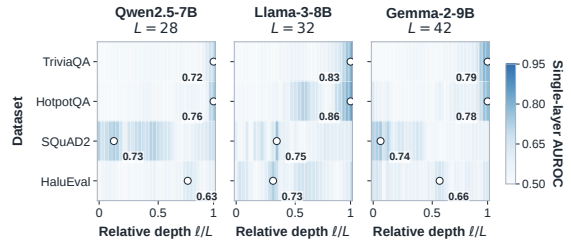


Figure 3: Single-layer AUROC heatmap from H_x^ℓ . Rows are benchmarks, columns are relative depth ℓ/L , and the peak cell in each row is circled. Peak-discriminative depth varies systematically across models.

et al., 2025). Appendix A.4 summarizes these baselines and the evaluation protocol used in our runs. All baseline results are produced under the same evaluation protocol, group-preserving 80/20 split, model suite, and model-dataset grid as TriLens.

4.2 Main Results

Table 1 reports test AUROC on all (model, dataset) cells against the six baselines. Our primary method, **TriLens (MLP)** (last row per model), uses the feature $\mathbf{f} = (H_a^\ell, H_m^\ell, H_x^\ell)_\ell$ with the same MLP probe topology used by ICR Probe. Unless otherwise noted, *TriLens* refers to this MLP variant used in Table 1.

Main finding. Hallucination risk is visible in the way internal certainty forms across layers: a compact trajectory of logit-lens entropies is sufficient to separate supported from hallucinated answers in a single forward pass. TriLens improves over ICR Probe on all 12 cells and is the highest-AUROC entry in 11 of the 12 cells, with gains ranging from +4.2 AUROC (Qwen2.5-7B on TriviaQA) to +16.0 AUROC (Qwen2.5-7B on SQuAD2). The improvement is consistent across model families (+10.2 averaged over Gemma-2, +11.0 over Qwen2.5, and +15.1 over Llama-3) and across benchmark regimes (+12.8 over the three context-grounded tasks and +10.2 on TriviaQA).

Probe and capacity controls. The pattern is not explained solely by the MLP probe. A linear probe on the same features also outperforms ICR Probe on all 12 cells, with the MLP adding another ~ 3.2 AUROC points (Appendix B). Dimension-matched controls further show that simply repeating H_x three times is nearly identical to H_x alone, whereas TriLens remains better on all 12 linear-probe cells. A stricter MLP sweep with L -dimensional reduc-

LLM	Methods	HaluEval	SQuAD2	HotpotQA	TriviaQA
Gemma-2-9B	PPL	0.5538	0.5348	0.7239	0.7536
	LN-Entropy	0.7357	0.6818	0.7349	0.7023
	LLM-Check	0.5780	0.5517	0.5102	0.5911
	SAPLMA	0.8123	0.7409	0.8329	0.7766
	SEP	0.6439	0.6627	0.6149	0.7726
	ICR Probe	<u>0.8436</u>	<u>0.8142</u>	<u>0.8409</u>	<u>0.8001</u>
	TriLens (Ours)	0.9277	0.9270	0.9433	0.9106
Qwen2.5-7B	PPL	0.5512	0.5278	0.6069	0.7041
	LN-Entropy	0.7286	0.6564	0.6913	0.6938
	LLM-Check	0.5367	0.5639	0.5518	0.5604
	SAPLMA	0.7725	0.7016	0.7689	0.8297
	SEP	0.6634	0.6418	0.6536	0.7449
	ICR Probe	<u>0.8076</u>	<u>0.7382</u>	<u>0.7865</u>	<u>0.7751</u>
	TriLens (Ours)	0.9053	0.9061	0.9242	<u>0.8103</u>
Llama-3-8B	PPL	0.5867	0.6472	0.6658	0.7085
	LN-Entropy	0.6574	0.6249	0.6661	0.5928
	LLM-Check	0.5263	0.5356	0.5559	0.5482
	SAPLMA	0.7315	0.7043	<u>0.7768</u>	0.7586
	SEP	0.7309	<u>0.7274</u>	0.6607	0.7048
	ICR Probe	<u>0.7671</u>	<u>0.7568</u>	0.7909	0.7396
	TriLens (Ours)	0.9136	0.9169	0.9425	0.8861

Table 1: **Main results:** test AUROC on four benchmark datasets under a matched evaluation protocol. **Bold** marks the best method in each model–dataset cell; underlining marks the second-best. TriLens improves over ICR Probe on all 12 cells and is best in 11, with average gain +12.1 AUROC. Per-cell standard deviations are listed in Appendix D.

tions of the three-pathway feature gives the same conclusion, indicating that the gain is not explained by input dimensionality alone.

4.3 Feature Ablation

The ablation answers whether the three readouts provide redundant or complementary evidence. We train the MLP probe on progressively larger subsets of \mathbf{f} and additionally include an intra-layer module-disagreement term $\text{JSD}_{am}^\ell = \text{JSD}(\text{lens}(a_i^\ell), \text{lens}(m_i^\ell))$ as a fifth feature, motivated by an intra-layer analog of ICR Probe’s consistency construction. Detailed ablation numbers are reported in Appendix E. The useful signal is not confined to the residual stream: adding either H_a or H_m to H_x improves performance on all 12 cells, and the full 3-entropy feature improves further in every cell. This indicates complementary predictive evidence across the three locations. Adding the intra-layer disagreement term changes AUROC by at most +0.006 per cell and +0.001 on average, suggesting little value beyond the entropy feature itself. The relative importance of H_a and H_m is architecture-dependent: the $H_a + H_x$ branch is stronger on Gemma-2, whereas $H_m + H_x$ is stronger on Llama-3 and part of Qwen2.5.

Figure 4 visualizes this complementarity from two angles. The left panel shows that the three

entropies open hallucination–correct separation at different depths and with architecture-specific pathway asymmetry; the right panel shows that both two-feature branches improve over H_x alone, while the full 3-entropy feature attains the highest overall mean.

4.4 Per-Layer Analysis

Detailed trajectory plots and peak-layer densities are deferred to Appendix F. In the main text, we focus on the summary heatmap in Figure 3. It shows that the peak-discriminative depth varies substantially across model–dataset pairs: some open-book cells peak in early or middle layers, while TriviaQA and several HotpotQA cells peak only near the final layer. This benchmark-conditional variation is consistent with prior mechanistic analyses (Geva et al., 2023; Zhang et al., 2025), which found that different models integrate knowledge at different depths.

4.5 Cross-Dataset Diagnostics

The main results in Table 1 use a separate probe per benchmark. We next use two benchmark-shift diagnostics; detailed results appear in Appendix C. First, we train one probe per model on the union of all four benchmark training splits and evaluate it on each benchmark’s held-out test split. This

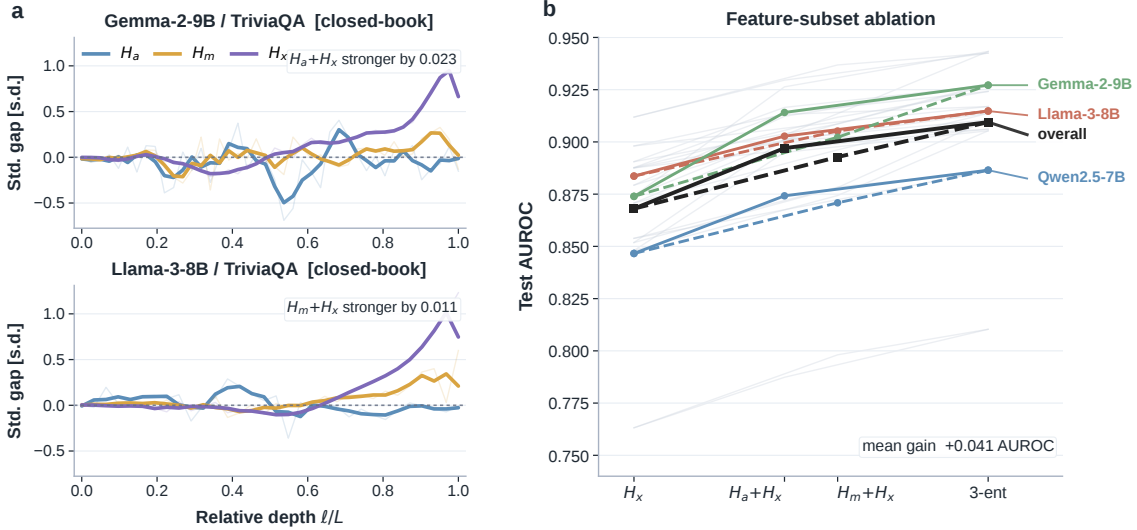


Figure 4: **Complementarity of the three entropy features.** **Left:** standardized per-layer separation between hallucinated and correct samples for H_a^ℓ , H_m^ℓ , and H_x^ℓ on two representative TriviaQA cells. The stronger two-feature branch is architecture-dependent. **Right:** feature-subset ablation across all 12 (model, dataset) cells. Colored lines show per-model means; grey lines show individual cells. Both two-feature branches improve over H_x alone, and the full three-entropy feature attains the best mean.

shared-probe setting still exceeds per-dataset ICR Probe on all 12 cells, with average gain +10.7 AUROC, while trailing our own per-dataset probes by only ~ 1.4 AUROC on average. Second, we train on one benchmark and test on another to form benchmark-level transfer heatmaps. In this cross-dataset transfer diagnostic, TriLens is the strongest method on all 12 off-diagonal cells and on all 16 cells overall; its off-diagonal mean AUROC is 0.848, compared with 0.738 for ICR Probe and 0.678 for SAPLMA. We therefore view the transfer result as further evidence that the gain is not confined to in-domain fitting, while still treating it as a distribution-shift diagnostic rather than as evidence of universal cross-benchmark generalization.

4.6 Comparison with Cross-Layer Contrast Signals

TriLens and DoLa (Chuang et al., 2024) both derive signals from logit-lens distributions but differ in what they measure: DoLa contrasts logit-lens distributions *across* layers and uses the result to modify decoding, while we measure entropy *within* each layer and use it for detection. To assess whether our per-layer entropy feature could be explained as a reformulation of the cross-layer contrast signal, we repurpose DoLa’s quantity as a detection feature and compare it head to head against ours.

Concretely, define the DoLa-style detection fea-

Dataset	DoLa-JSD	TriLens	TriLens+DoLa
HaluEval	0.7623	0.9136	0.9155
SQuAD2	0.8538	0.9158	0.9172
HotpotQA	0.8928	0.9425	0.9422
TriviaQA	0.8679	0.8861	0.8860

Table 2: Head-to-head comparison against DoLa-style cross-layer contrast features (MLP probe, test AUROC). Full results appear in Appendix G.

ture as

$$f_{\text{DoLa}}^\ell = \text{JSD}(\text{lens}(x^\ell), \text{lens}(x^L)), \quad (4)$$

$$\ell = 0, \dots, L - 1.$$

and assemble the L -dimensional vector \mathbf{f}_{DoLa} . We train the same MLP probe on (a) \mathbf{f}_{DoLa} , (b) the TriLens feature $\mathbf{f}_{3\text{-ent}}$, and (c) their concatenation $\mathbf{f}_{3\text{-ent}} \oplus \mathbf{f}_{\text{DoLa}}$. Experiments are run under the same group-preserving 80/20 evaluation protocol used in our main experiments, using the same 10k-scale feature files as the main paper. Table 2 reports the main-setting comparison, while Appendix G gives the full configuration grid.

Findings. Table 2 presents three observations. First, **our per-layer entropy consistently outperforms the DoLa-style cross-layer contrast feature on every dataset**, with gaps ranging from +0.02 to +0.15 AUROC and an average gain of roughly +0.07 in the main setting. Second, **their**

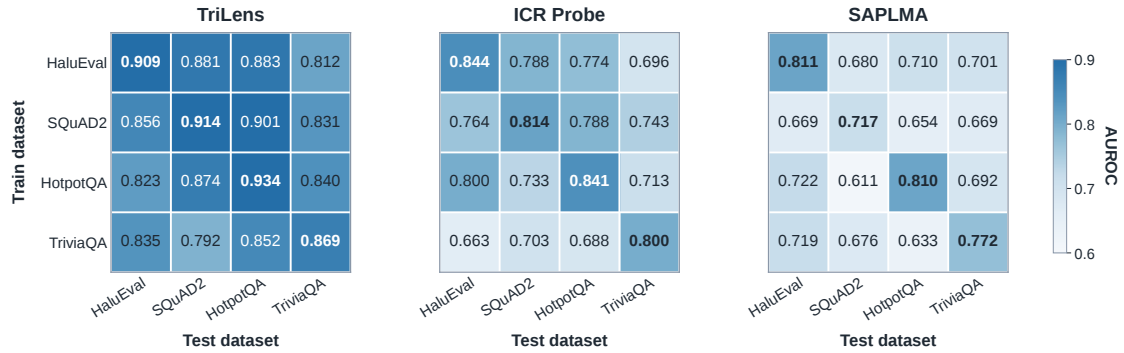


Figure 5: Cross-dataset generalization heatmaps for TriLens, ICR Probe, and SAPLMA. Each cell reports AUROC for the train-on-row, test-on-column setting. TriLens is strongest on all 12 off-diagonal cells and on all 16 cells overall, indicating more stable transfer under benchmark shift.

union changes performance only slightly and without a consistent upward pattern relative to our feature alone: averaged across the four datasets, the union changes AUROC by less than +0.001. The cross-layer contrast signal therefore contributes little beyond what is already captured by per-layer entropy. Third, the DoLa-style feature degrades most severely on HaluEval, where in our default setup, intermediate layers have already converged onto the final layer’s prediction (so the JSD collapses toward zero regardless of whether the response is hallucinated or correct); per-layer entropy is not subject to this degeneracy because it measures the shape of each layer’s distribution independently.

5 Conclusion

We propose *TriLens*, a 3L-dimensional feature that computes Shannon entropies of logit-lens distributions at the MHSA output, FFN output, and residual stream at every layer. Across 3 LLMs \times 4 benchmarks, under the same MLP probe family used by a strong prior baseline, this feature improves over ICR Probe on all 12 cells with an average gain of +12.1 AUROC and attains the highest AUROC in 11 of the 12 cells. Ablation identifies all three module entropies as independently informative and rules out intra-layer module-disagreement as redundant. A head-to-head comparison suggests that our feature captures much of the useful signal in DoLa-style cross-layer contrast. Layer-wise analysis further indicates that the most informative depth varies across architectures and benchmarks rather than concentrating at a universal layer. Taken together, these findings suggest that simple, layer-wise features merit consideration alongside more

elaborate probe architectures in white-box hallucination detection.

Limitations

TriLens provides an effective white-box detection signal, but it has several limitations. First, it requires access to internal activations and is therefore restricted to open-source or otherwise transparent models. Second, our study focuses on detection rather than mitigation, so it does not directly address how such signals should be used to reduce hallucinations during generation. Finally, although our experiments span multiple benchmarks and model families, extending the evaluation to broader model scales and more complex generation settings would provide a stronger test of robustness and generality.

Ethics Statement

Our study is conducted entirely on publicly available LLMs and established public benchmarks, introduces no new human-subject data collection, and requires no additional annotation. We do not release any new text dataset containing personally identifying information or offensive content; all reported results are aggregate metrics over benchmark splits. We use the public models and benchmarks only under their publicly stated access conditions and licenses or terms of use, and our released code contains no redistributed model weights or benchmark text. Our use of these artifacts is limited to research evaluation of hallucination detection and is consistent with their benchmark or model-evaluation purpose and access conditions. We acknowledge the standard dual-use consideration that more accurate hallucination detection could be used to identify, and thereby refine, adversarial gen-

566	erations. However, as a detection-only capability	Dominican Republic. Association for Computational	619
567	without any generation or mitigation component,	Linguistics.	620
568	the net effect is defensive.		
	References		
569			
570	Amos Azaria and Tom Mitchell. 2023. The internal	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	621
571	state of an LLM knows when it’s lying . In <i>Find-</i>	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	622
572	<i>ings of the Association for Computational Linguistics:</i>	Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,	623
573	<i>EMNLP 2023</i> , pages 967–976, Singapore. Associa-	Alex Vaughan, and 1 others. 2024. The llama 3 herd	624
574	tion for Computational Linguistics.	of models. <i>arXiv preprint arXiv:2407.21783</i> .	625
575	Nora Belrose, Igor Ostrovsky, Lev McKinney, Zach Fur-	Xiaomeng Hu, Yiming Zhang, Ru Peng, Haozhe Zhang,	626
576	man, Logan Smith, Danny Halawi, Stella Biderman,	Chenwei Wu, Gang Chen, and Junbo Zhao. 2024.	627
577	and Jacob Steinhardt. 2023. Eliciting latent predic-	Embedding and gradient say wrong: A white-box	628
578	tions from transformers with the tuned lens. <i>arXiv</i>	method for hallucination detection . In <i>Proceedings</i>	629
579	<i>preprint arXiv:2303.08112</i> .	<i>of the 2024 Conference on Empirical Methods in</i>	630
580	Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and	<i>Natural Language Processing</i> , pages 1950–1959, Mi-	631
581	Chris Kedzie. 2024. Do androids know they’re only	ami, Florida, USA. Association for Computational	632
582	dreaming of electric sheep? In <i>Findings of the As-</i>	Linguistics.	633
583	<i>sociation for Computational Linguistics: ACL 2024</i> ,	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	634
584	pages 4401–4420, Bangkok, Thailand. Association	Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	635
585	for Computational Linguistics.	Madotto, and Pascale Fung. 2023. Survey of hal-	636
586	Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu,	lucination in natural language generation. <i>ACM com-</i>	637
587	Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024.	<i>puting surveys</i> , 55(12):1–38.	638
588	Inside: LLMs’ internal states retain the power of hal-	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke	639
589	lucination detection. In <i>The Twelfth International</i>	Zettlemoyer. 2017. TriviaQA: A large scale distantly	640
590	<i>Conference on Learning Representations</i> .	supervised challenge dataset for reading comprehen-	641
591	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon	sion . In <i>Proceedings of the 55th Annual Meeting of</i>	642
592	Kim, James R. Glass, and Pengcheng He. 2024. Dola:	<i>the Association for Computational Linguistics (Vol-</i>	643
593	Decoding by contrasting layers improves factuality in	<i>ume 1: Long Papers)</i> , pages 1601–1611, Vancouver,	644
594	large language models . In <i>The Twelfth International</i>	Canada. Association for Computational Linguistics.	645
595	<i>Conference on Learning Representations</i> .	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	646
596	Kamil Ciosek, Nicolò Felicioni, and Sina Ghiassian.	Henighan, Dawn Drain, Ethan Perez, Nicholas	647
597	2025. Hallucination detection on a budget: Effi-	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	648
598	cient bayesian estimation of semantic entropy. <i>arXiv</i>	Tran-Johnson, and 1 others. 2022. Language mod-	649
599	<i>preprint arXiv:2504.03579</i> .	els (mostly) know what they know. <i>arXiv preprint</i>	650
600	Nelson Elhage, Neel Nanda, Catherine Olsson, and 1	<i>arXiv:2207.05221</i> .	651
601	others. 2021. A mathematical framework for trans-	Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa	652
602	former circuits . <i>Anthropic Technical Report</i> .	Schut, Shreshth A. Malik, and Yarin Gal. 2024. Se-	653
603	Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and	mantic entropy probes: Robust and cheap hallucina-	654
604	Yarin Gal. 2024. Detecting hallucinations in large	tion detection in llms . <i>ArXiv</i> , abs/2406.15927.	655
605	language models using semantic entropy . <i>Nature</i> ,	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	656
606	630:625 – 630.	Semantic uncertainty: Linguistic invariances for un-	657
607	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir	certainty estimation in natural language generation.	658
608	Globerson. 2023. Dissecting recall of factual associa-	<i>arXiv preprint arXiv:2302.09664</i> .	659
609	tions in auto-regressive language models . In <i>Proceed-</i>	Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and	660
610	<i>ings of the 2023 Conference on Empirical Methods in</i>	Ji-Rong Wen. 2023. HaluEval: A large-scale hal-	661
611	<i>Natural Language Processing</i> , pages 12216–12235,	lucination evaluation benchmark for large language	662
612	Singapore. Association for Computational Linguis-	models . In <i>Proceedings of the 2023 Conference on</i>	663
613	tics.	<i>Empirical Methods in Natural Language Processing</i> ,	664
614	Mor Geva, Roei Schuster, Jonathan Berant, and Omer	pages 6449–6464, Singapore. Association for Com-	665
615	Levy. 2021. Transformer feed-forward layers are key-	putational Linguistics.	666
616	value memories . In <i>Proceedings of the 2021 Confer-</i>	Andrey Malinin and Mark Gales. 2020. Uncertainty esti-	667
617	<i>ence on Empirical Methods in Natural Language Pro-</i>	mation in autoregressive structured prediction. <i>arXiv</i>	668
618	<i>cessing</i> , pages 5484–5495, Online and Punta Cana,	<i>preprint arXiv:2002.07650</i> .	669
		Dang Nguyen, Ali Payani, and Baharan Mirzasoleiman.	670
		2025. Beyond semantic entropy: Boosting llm uncer-	671
		tainty quantification with pairwise semantic similar-	672
		ity. In <i>Findings of the Association for Computational</i>	673
		<i>Linguistics: ACL 2025</i> , pages 4530–4540.	674

675	nostalgebraist. 2020. Interpreting GPT: The	In <i>Proceedings of the 2018 Conference on Empirical</i>	732
676	logit lens. https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens .	<i>Methods in Natural Language Processing</i> , pages	733
677		2369–2380, Brussels, Belgium. Association for Com-	734
678		putational Linguistics.	735
679	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.	Zhenliang Zhang, Xinyu Hu, Huixuan Zhang, Junzhe	736
680	Know what you don’t know: Unanswerable ques-	Zhang, and Xiaojun Wan. 2025. ICR probe: Track-	737
681	tions for SQuAD. In <i>Proceedings of the 56th Annual</i>	ing hidden state dynamics for reliable hallucination	738
682	<i>Meeting of the Association for Computational Lin-</i>	detection in LLMs. In <i>Proceedings of the 63rd An-</i>	739
683	<i>guistics (Volume 2: Short Papers)</i> , pages 784–789,	<i>annual Meeting of the Association for Computational</i>	740
684	Melbourne, Australia. Association for Computational	<i>Linguistics (ACL)</i> .	741
685	Linguistics.		
686	Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mo-	A Experimental Details	742
687	hammad Saleh, Balaji Lakshminarayanan, and Pe-	A.1 Computing Infrastructure	743
688	ter J Liu. 2022. Out-of-distribution detection and	Our experiments were conducted on a server	744
689	selective generation for conditional language models.	equipped with 6 NVIDIA GeForce RTX 4090	745
690	<i>arXiv preprint arXiv:2209.15558</i> .	GPUs (24 GB memory each), running PyTorch	746
691	Albert Sawczyn, Jakub Binkowski, Denis Janiak, Bog-	with CUDA acceleration. Across multiple rounds	747
692	dan Gabrys, and Tomasz Jan Kajdanowicz. 2026.	of feature extraction, ablation, and evaluation exper-	748
693	FactSelfCheck: Fact-level black-box hallucination	iments, the total computational budget amounted	749
694	detection for LLMs. In <i>Findings of the Association</i>	to 100–200 GPU hours. Feature extraction re-	750
695	<i>for Computational Linguistics: EACL 2026</i> , pages	quires one teacher-forced forward pass per scored	751
696	5603–5621, Rabat, Morocco. Association for Com-	sequence; probe training is lightweight relative to	752
697	putational Linguistics.	LLM-state extraction and is repeated across seeds	753
698	Ahmad Shapiro, Karan Taneja, and Ashok Goel. 2026.	and evaluation configurations reported in the main	754
699	Halt: Hallucination assessment via log-probs as time	text and appendix.	755
700	series. <i>arXiv preprint arXiv:2602.02888</i> .		
701	Gaurang Sriramanan, Siddhant Bharti, Vinu Sankar	A.2 Details about Datasets	756
702	Sadasivan, Shoumik Saha, Priyatham Kattakinda,	For our experiments, we randomly sampled 10,000	757
703	and Soheil Feizi. 2024. Llm-check: Investigating	instances from each dataset. Specifically, we used	758
704	detection of hallucinations in large language mod-	the QA subset of HaluEval and the ‘rc.nocontext’	759
705	els. In <i>Advances in Neural Information Processing</i>	subset of TriviaQA. These benchmarks are English-	760
706	<i>Systems</i> , volume 37, pages 34188–34216. Curran As-	language QA-style evaluation sets covering factoid	761
707	sociates, Inc.	question answering, unanswerable questions, multi-	762
708	Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya	hop reasoning, and hallucinated QA responses; we	763
709	Sachan. 2023. A mechanistic interpretation of arith-	do not use demographic annotations or perform de-	764
710	metic reasoning in language models using causal	demographic subgroup analysis. Each dataset is split	765
711	mediation analysis. In <i>Proceedings of the 2023 Con-</i>	into 80%-20% for training and testing. We train	766
712	<i>ference on Empirical Methods in Natural Language</i>	and test on each dataset, reporting the correspond-	767
713	<i>Processing</i> , pages 7035–7052, Singapore. Associa-	ing results. This experimental setup for baseline	768
714	tion for Computational Linguistics.	methods matches ours exactly.	769
715	Gemma Team, Morgane Riviere, Shreya Pathak,	A.3 Sampling, Splits, and Supervision	770
716	Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-	All trainable methods use group-preserving 80/20	771
717	raju, Léonard Hussenot, Thomas Mesnard, Bobak	train/test splits, and the same split assignment is	772
718	Shahriari, Alexandre Ramé, and 1 others. 2024.	shared by TriLens and all baselines. All reported	773
719	Gemma 2: Improving open language models at a	results are averaged over five random seeds.	774
720	practical size. <i>arXiv preprint arXiv:2408.00118</i> .		
721	Qwen An Yang, Baosong Yang, Beichen Zhang,	A.4 Details about Baselines	775
722	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan	We summarize the six baselines in Table 1. The	776
723	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	first three are training-free sample-level scorers,	777
724	ran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei	whereas the latter three are trainable white-box	778
725	Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Jun-	detectors. All are evaluated in the same pipeline as	779
726	yang Lin, and 25 others. 2024. Qwen2.5 technical		
727	report. <i>ArXiv</i> , abs/2412.15115.		
728	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,		
729	William Cohen, Ruslan Salakhutdinov, and Christo-		
730	pher D. Manning. 2018. HotpotQA: A dataset for		
731	diverse, explainable multi-hop question answering.		

TriLens, on the same model–dataset grid and group-preserving 80/20 split. Scalar baselines are used directly as sample-level scores; trainable baselines retain their original feature definition and are fit on the corresponding training split. For the 10k-scale evaluation sets used in the main paper, HaluEval and SQuAD2 follow the dataset sizes reported in the main text.

PPL. Perplexity of the scored sequence under the same forward-pass setup. Higher perplexity indicates lower model confidence and is used directly as a hallucination score.

LN-Entropy. Length-normalized predictive entropy of the scored sequence. This baseline measures output uncertainty while partially controlling for raw response length.

LLM-Check. A training-free attention-kernel baseline that scores responses from structural statistics of self-attention kernels rather than hidden states.

SAPLMA. A supervised hidden-state baseline that trains an MLP classifier on internal activations from selected layers.

SEP. A probe-based hidden-state baseline motivated by semantic entropy, using learned detectors over internal representations rather than black-box semantic clustering over sampled outputs.

ICR Probe. The strongest prior white-box baseline. It derives one scalar feature per layer from residual-stream update statistics and feeds the resulting layer-wise vector to the same 4-layer MLP architecture used throughout our comparisons; the shared probe hyperparameters are listed in Appendix H.

A.5 Probe Comparability Note

TriLens and ICR Probe are compared under the same training protocol and the same hidden-layer MLP topology, but not under equal input dimensionality: TriLens uses $3L$ input features whereas ICR Probe uses L . The main text therefore interprets the comparison as a feature-plus-probe comparison within a matched probe family, rather than as a fully dimension-controlled capacity study.

B Probe Controls and Fairness Checks

Table 3 reports the same TriLens feature used in the main paper, replacing the MLP probe with L2-

regularized logistic regression under the same feature configuration as the main results. The qualitative conclusion is unchanged: the linear probe already improves over ICR Probe on all 12 cells, with the MLP adding a further ~ 3.2 AUROC points on average. To isolate the effect of input dimensionality, Table 4 adds a dimension-matched control $H_x^{\times 3}$ that repeats the residual-stream entropy feature three times, yielding the same $3L$ input size as TriLens without introducing any new signal. Table 5 extends this check under the same MLP probe family with stricter dimensionality controls: three handcrafted L -dimensional reductions of the three-pathway feature and a train-split PCA projection of the full $3L$ -dimensional TriLens vector back to L . We further report alternative readout controls in Appendix B.1, comparing the pathway-separated TriLens readout against intermediate-state and additive-write readouts built from the same layer-wise logit-lens pipeline.

B.1 Alternative Readout Controls

The main TriLens feature applies the final logit-lens readout to the isolated MHSA write a^ℓ , the isolated FFN write m^ℓ , and the composed residual state x^ℓ . To test whether the gains depend on this particular choice of readout location, we compare it with two more conservative alternatives built from the same layer-wise pipeline: (i) H_{pre} , the entropy of the intermediate state $x^{\ell-1} + a^\ell$ after the attention write but before the FFN write, and (ii) H_p , the entropy of the additive-write state $a^\ell + m^\ell$. We then pair each alternative with the standard residual-stream readout H_x and evaluate the resulting features under the same protocol as the main paper.

Three patterns are consistent across probe families. First, performance improves monotonically as the readout moves from residual-only H_x to $H_{\text{pre}} + H_x$, then to $H_p + H_x$, and finally to the full pathway-separated TriLens feature. Under the linear probe, the corresponding average AUROC increases from 0.814 to 0.833, 0.860, and 0.878, with TriLens strongest on all 12 cells. Second, the same ordering largely persists under the MLP probe, where TriLens is strongest on 11 of 12 cells and improves over H_x by +0.052 AUROC on average. Third, the mixed-state alternatives do help, but they remain consistently below the explicit three-way decomposition. We therefore interpret this sweep as evidence that TriLens benefits from separating MHSA, FFN, and residual-stream readouts, rather than from using an arbitrary auxiliary lens

LLM	Methods	HaluEval	SQuAD2	HotpotQA	TriviaQA
Gemma-2-9B	ICR Probe	0.8436	0.8142	0.8409	0.8001
	TriLens (Linear)	0.8963	0.9004	0.9203	0.8864
Qwen2.5-7B	ICR Probe	0.8003	0.7456	0.7917	0.7684
	TriLens (Linear)	0.8597	0.8696	0.8746	0.7780
Llama-3-8B	ICR Probe	0.7603	0.7634	0.7982	0.7325
	TriLens (Linear)	0.8731	0.8806	0.9253	0.8700

Table 3: Linear-probe results on the same 12 cells as Table 1. All TriLens entries are 5-seed means; seed standard deviation is ≤ 0.007 in every cell.

LLM	Dataset	H_x	$H_x^{\times 3}$	TriLens
Gemma-2-9B	HaluEval	0.7970	0.7970	0.8963
	SQuAD2	0.8497	0.8497	0.9004
	HotpotQA	0.7812	0.7813	0.9203
	TriviaQA	0.8238	0.8238	0.8864
Qwen2.5-7B	HaluEval	0.7746	0.7746	0.8597
	SQuAD2	0.8389	0.8389	0.8696
	HotpotQA	0.7992	0.7992	0.8746
	TriviaQA	0.7384	0.7384	0.7780
Llama-3-8B	HaluEval	0.8247	0.8248	0.8731
	SQuAD2	0.8254	0.8255	0.8806
	HotpotQA	0.8717	0.8717	0.9253
	TriviaQA	0.8385	0.8385	0.8700

Table 4: Dimension-matched fairness control under the linear probe. $H_x^{\times 3}$ repeats the L -dimensional residual-stream entropy feature three times to match TriLens’s $3L$ input size. The repeated baseline is nearly identical to H_x , while TriLens remains better on all 12 cells.

location.

C Cross-Dataset Generalization Heatmaps

The main paper briefly reports the shared-probe variant of TriLens. Table 7 gives the corresponding per-benchmark numbers. We additionally report benchmark-level cross-dataset generalization heatmaps for TriLens, ICR Probe, and SAPLMA. In each 4×4 matrix, a probe is trained on the row benchmark and evaluated on the column benchmark under the same train/test transfer protocol.

Aggregate statistics. Across the 12 off-diagonal cells, TriLens reaches a mean AUROC of 0.848, compared with 0.738 for ICR Probe and 0.678 for SAPLMA. Its diagonal mean is also higher (0.907 vs. 0.825 and 0.778), so the advantage is not limited to either in-domain or out-of-domain evaluation alone.

Transfer pattern. The transfer gap is not perfectly uniform across source benchmarks. For

TriLens, probes trained on SQuAD2 or HaluEval give the strongest off-diagonal averages (0.863 and 0.859), while TriviaQA is the weakest source row (0.826 off-diagonal mean). Even so, all 12 off-diagonal TriLens cells remain above 0.79, whereas the baselines show broader degradation. We therefore interpret Figure 5 as evidence of comparatively stable transfer under benchmark shift, while still viewing cross-dataset evaluation as a diagnostic rather than as a substitute for the in-domain results in Table 1.

D Main-Table Standard Deviations

Table 8 lists the standard deviations for the TriLens entries in Table 1, across the same five seeds used in the main-paper evaluation.

E Feature Ablation

Table 9 reports the full feature-set ablation from the main paper discussion.

Table 10 lists the standard deviations across the same five seeds used in Table 9.

F Per-Layer Diagnostics and Mechanistic Reading

This appendix collects the layer-wise plots and the fuller interpretation omitted from the compressed main text.

Peak-layer analysis is included only as a descriptive diagnostic rather than as a tuned component of TriLens. The main results always use the full per-layer feature vector, without selecting or reweighting layers based on validation performance. We report the peak-discriminative layer only to summarize where the single-feature H_x^ℓ detector reaches its highest AUROC in each model-dataset cell.

Context-grounded vs. closed-book regimes. On context-grounded benchmarks such as SQuAD2 and HotpotQA, the correct-hallucinated gap often

LLM	Dataset	H_x	$H_x^{\times 3}$	TriMeanL	TriMaxL	TriMinL	TriLens+PCA $\rightarrow L$	TriLens
Gemma-2-9B	HaluEval	0.8325	0.8420	0.8364	0.8140	0.8188	0.8794	0.8967
	SQuAD2	0.8532	0.8638	0.8461	0.8319	0.8349	0.8773	0.8951
	HotpotQA	0.8051	0.8272	0.8326	0.7991	0.8050	0.8927	0.9145
	TriviaQA	0.8237	0.8313	0.8057	0.7681	0.8053	0.8714	0.8838
Qwen2.5-7B	HaluEval	0.8000	0.8131	0.8004	0.7534	0.7877	0.8379	0.8576
	SQuAD2	0.8523	0.8606	0.8358	0.8134	0.8202	0.8669	0.8743
	HotpotQA	0.8376	0.8529	0.8155	0.7813	0.7585	0.8655	0.8789
	TriviaQA	0.7330	0.7387	0.6933	0.6218	0.6934	0.7667	0.7700
Llama-3-8B	HaluEval	0.7635	0.7783	0.7874	0.7480	0.7840	0.8209	0.8379
	SQuAD2	0.8100	0.8383	0.7902	0.7469	0.8077	0.8093	0.8448
	HotpotQA	0.8612	0.8710	0.8737	0.8057	0.8539	0.9034	0.9156
	TriviaQA	0.8305	0.8326	0.8130	0.6753	0.8110	0.8471	0.8578

Table 5: Stricter dimensionality controls under the MLP probe. Besides $H_x^{\times 3}$, we compare three handcrafted L -dimensional reductions of the three-pathway feature and a train-split PCA projection of the full $3L$ -dimensional TriLens vector back to L . Full TriLens remains strongest on all 12 cells.

Readout	Linear		MLP	
	Avg	Wins	Avg	Wins
H_x	0.8137	0/12	0.8169	0/12
$H_{\text{pre}} + H_x$	0.8334	0/12	0.8268	0/12
$H_p + H_x$	0.8604	0/12	0.8496	1/12
TriLens	0.8779	12/12	0.8690	11/12

Table 6: Readout-location robustness summary across the 12 model–dataset cells. “Avg” denotes mean test AU-ROC; “Wins” counts how often a readout is the strongest cell-wise configuration. All entries are 5-seed means under the same protocol as the main paper.

Model	Test	ICR [†]	Multi (ours)	Δ
Qwen2.5-7B	HaluEval	0.8003	0.8856	+0.085
	SQuAD2	0.7456	0.8915	+0.146
	HotpotQA	0.7917	0.9102	+0.118
	TriviaQA	0.7684	0.7998	+0.031
Llama-3-8B	HaluEval	0.7603	0.8965	+0.136
	SQuAD2	0.7634	0.9045	+0.141
	HotpotQA	0.7982	0.9329	+0.135
	TriviaQA	0.7325	0.8739	+0.141
Gemma-2-9B	HaluEval	0.8436	0.9132	+0.070
	SQuAD2	0.8142	0.9129	+0.099
	HotpotQA	0.8409	0.9225	+0.082
	TriviaQA	0.8001	0.8971	+0.097
avg		0.7883	0.8950	+0.107

Table 7: **Multi-dataset probe**: one probe per model, trained on the union of all four benchmarks and evaluated on each held-out test split. Each cell is a 5-seed mean (std ≤ 0.008). [†]ICR entries are our per-dataset reruns from Table 1.

933 opens earlier and remains visible across much of
934 the upper stack. On TriviaQA, by contrast, the tra-
935 jectories tend to remain close until the final layers,
936 consistent with a closed-book setting in which the
937 relevant parametric recall signal is resolved late.

938 **Architecture-dependent peak depth.** Figure 3
939 shows that the peak-discriminative layer is not uni-
940 versal across architectures. Some Qwen2.5 cells
941 peak in earlier or mid layers, whereas Llama-3 and
942 Gemma-2 frequently peak closer to the final layers.
943 This pattern is consistent with prior observations
944 that different model families accumulate factual
945 recall at different depths.

946 **Pathway asymmetry.** The reversal between
947 Gemma-2 and Llama-3 in Figure 4 suggests that
948 the relative usefulness of H_a and H_m depends on
949 how readable each pathway’s contribution is under
950 the model’s own unembedding. Gemma-2’s hy-
951 brid attention stack may shift that balance toward
952 MHSA relative to the other two models.

953 **Smallest-gain cell.** The smallest gain in the main
954 table occurs on Qwen2.5-7B / TriviaQA. Its late-
955 layer trajectories show weaker separation than the
956 corresponding cells for Llama-3 and Gemma-2,
957 suggesting a more diffuse parametric recall signal
958 in this closed-book setting.

959 Table 11 reports the peak-discriminative layer
960 for the single-feature detector based on H_x^ℓ , using
961 the same setup as Figures 6–3. Layer indices are
962 zero-based to match the figure annotations.

963 G DoLa Defense: Full Configuration 964 Grid

965 Table 12 provides the full configuration grid behind
966 Table 2: both probe families and all four datasets,
967 using the same feature configuration as the main

Model	Dataset	Std
Gemma-2-9B	HaluEval	0.0022
Gemma-2-9B	SQuAD2	0.0025
Gemma-2-9B	HotpotQA	0.0032
Gemma-2-9B	TriviaQA	0.0027
Qwen2.5-7B	HaluEval	0.0053
Qwen2.5-7B	SQuAD2	0.0025
Qwen2.5-7B	HotpotQA	0.0015
Qwen2.5-7B	TriviaQA	0.0051
Llama-3-8B	HaluEval	0.0014
Llama-3-8B	SQuAD2	0.0037
Llama-3-8B	HotpotQA	0.0033
Llama-3-8B	TriviaQA	0.0037

Table 8: Seed standard deviations for the TriLens entries in Table 1.

Model	Data	H_x	H_a+H_x	H_m+H_x	3-ent	3-ent+JSD $_{am}^{\ell}$
Gemma-2	HaluEval	.8793	.9123	.9087	.9277	.9286
	SQuAD2	.8905	.9165	.9085	.9270	.9285
	HotpotQA	.8741	.9264	.9135	.9433	.9451
	TriviaQA	.8518	.9013	.8782	.9106	.9115
Qwen2.5	HaluEval	.8480	.8950	.8743	.9053	.9067
	SQuAD2	.8876	.9006	.8992	.9061	.9065
	HotpotQA	.8877	.9138	.9118	.9242	.9256
	TriviaQA	.7632	.7874	.7981	.8103	.8102
Llama-3	HaluEval	.8708	.9008	.8962	.9136	.9194
	SQuAD2	.8981	.9132	.9095	.9169	.9194
	HotpotQA	.9119	.9296	.9368	.9425	.9444
	TriviaQA	.8538	.8675	.8783	.8861	.8846

Table 9: Feature-set ablation (MLP probe, test AUROC). **3-ent** = (H_a, H_m, H_x) ; the last column adds JSD $_{am}^{\ell}$.

paper.

H Efficiency and Complexity Analysis

H.1 Training Hyperparameters

H.2 Time Complexity

For a decoder with L layers and hidden size d , TriLens adds a logit-lens readout at three locations per layer and reduces each readout to a scalar entropy. Relative to methods that store or classify high-dimensional hidden states, the resulting sample representation is compact: TriLens produces a $3L$ -dimensional vector, whereas hidden-state baselines typically operate on features that scale with d or selected layer subsets of size comparable to d . Probe-time complexity is therefore negligible compared with the upstream LLM forward pass; in practice, extraction dominates runtime and the downstream classifier is cheap to retrain across seeds.

Model	Data	H_x	H_a+H_x	H_m+H_x	3-ent	3-ent+JSD $_{am}^{\ell}$
Gemma-2	HaluEval	0.004	0.002	0.003	0.002	0.002
	SQuAD2	0.006	0.004	0.006	0.004	0.004
	HotpotQA	0.003	0.003	0.003	0.004	0.004
	TriviaQA	0.004	0.004	0.004	0.004	0.005
Qwen2.5	HaluEval	0.005	0.002	0.005	0.003	0.003
	SQuAD2	0.006	0.005	0.004	0.006	0.005
	HotpotQA	0.005	0.005	0.003	0.004	0.005
	TriviaQA	0.006	0.004	0.008	0.006	0.006
Llama-3	HaluEval	0.003	0.003	0.001	0.002	0.002
	SQuAD2	0.003	0.004	0.004	0.004	0.004
	HotpotQA	0.004	0.003	0.004	0.002	0.003
	TriviaQA	0.003	0.003	0.005	0.004	0.004

Table 10: Standard deviations corresponding to Table 9 (MLP probe, same aggregation as the main results).

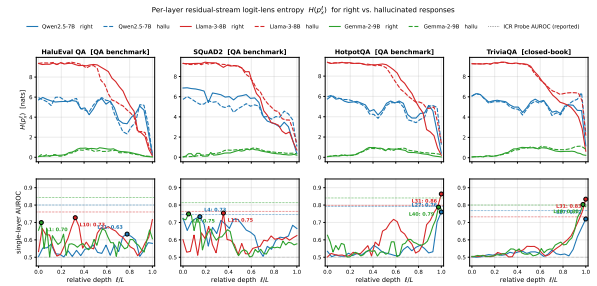


Figure 6: Per-layer residual-stream logit-lens entropy H_x^{ℓ} (top) and single-layer AUROC (bottom) across three models and four benchmarks. Solid lines denote correct responses; dashed lines denote hallucinated responses.

H.3 Memory Footprint of Logit-Lens Extraction

For a model with vocabulary size $|V|$, a naive implementation could materialize, for each scored token, $3L$ intermediate logit-lens distributions of the form $\text{softmax}(W_U \cdot \text{Norm}(z)) \in \mathbb{R}^{|V|}$. In float32, this corresponds to a peak memory cost of $3L|V| \times 4$ bytes per token if all such distributions are retained simultaneously. For example, for Llama-3-8B ($L=32$, $|V| \approx 128\text{K}$), the raw peak is approximately $3 \times 32 \times 128000 \times 4 \approx 49$ MB per token. In our implementation, however, each distribution is used immediately to compute its Shannon entropy and is then discarded, so the persistent storage scales as $O(3L)$ rather than $O(3L|V|)$. For Llama-3-8B, this reduces the retained per-token state to only 96 scalar entropy values.

Model	Dataset	ℓ^*	AUROC
Qwen2.5-7B	HaluEval	L21	0.6332
Qwen2.5-7B	SQuAD2	L4	0.7329
Qwen2.5-7B	HotpotQA	L27	0.7603
Qwen2.5-7B	TriviaQA	L27	0.7199
Llama-3-8B	HaluEval	L10	0.7274
Llama-3-8B	SQuAD2	L11	0.7541
Llama-3-8B	HotpotQA	L31	0.8638
Llama-3-8B	TriviaQA	L31	0.8340
Gemma-2-9B	HaluEval	L1	0.6994
Gemma-2-9B	SQuAD2	L2	0.7491
Gemma-2-9B	HotpotQA	L40	0.7881
Gemma-2-9B	TriviaQA	L40	0.8033

Table 11: Peak-discriminative layer ℓ^* and the corresponding single-layer AUROC for H_x^ℓ .

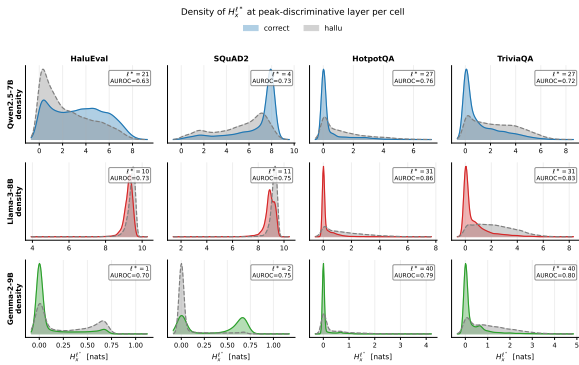


Figure 7: Probability density of $H_x^{\ell^*}$ at the peak-discriminative layer for correct vs. hallucinated responses.

Probe	Dataset	DoLa-JSD	TriLens	TriLens+DoLa
MLP	HaluEval	0.7623	0.9136	0.9155
	SQuAD2	0.8538	0.9158	0.9172
	HotpotQA	0.8928	0.9425	0.9422
	TriviaQA	0.8679	0.8861	0.8860
Linear	HaluEval	0.6790	0.8731	0.8738
	SQuAD2	0.7189	0.8778	0.8798
	HotpotQA	0.8111	0.9253	0.9256
	TriviaQA	0.8383	0.8700	0.8710

Table 12: Full DoLa comparison grid under both probe families.

Model weights dtype	bfloat16
Logit-lens projection dtype	float32
Attention implementation	SDPA
Temperature τ	1.0
Max response tokens	32
Train/test split	80/20 group-preserving split
Seeds	5
Main-paper dataset sizes	reported in Section 4.1
Row labels per dataset	exactly balanced correct/hallucinated labels
Linear probe	L2 logreg, $C = 1$
MLP probe architecture	$3L \rightarrow 128 \rightarrow 64 \rightarrow 32 \rightarrow 1$
MLP activation	LeakyReLU(0.01)
MLP regularization	BatchNorm + Dropout 0.3
MLP optimizer	Adam, lr = 5×10^{-4}
MLP training	50 epochs, batch 32
MLP LR schedule	ReduceLROnPlateau ($p=5, f=0.5$)

Table 13: Training hyperparameters used in the main experiments.