

LatentAudit: Real-Time White-Box Faithfulness Monitoring for Retrieval-Augmented Generation with Verifiable Deployment

Anonymous authors

Paper under double-blind review

Abstract

1 Retrieval-augmented generation (RAG) mitigates hallucination but does
2 not eliminate it: a deployed system must still decide, at inference time,
3 whether its answer is actually supported by the retrieved evidence. We
4 introduce *LatentAudit*, a white-box auditor that pools mid-to-late residual-
5 stream activations from an open-weight generator and measures their Ma-
6 halanobis distance to the evidence representation. The resulting quadratic
7 rule requires no auxiliary judge model, runs at generation time, and is
8 simple enough to calibrate on a small held-out set. We show that residual-
9 stream geometry carries a usable faithfulness signal, that this signal sur-
10 vives architecture changes and realistic retrieval failures, and that the same
11 rule remains amenable to public verification. On PubMedQA with Llama-3-
12 8B, LatentAudit reaches 0.942 AUROC with 0.77 ms overhead. Across three
13 QA benchmarks and five model families (Llama-2/3, Qwen-2.5/3, Mistral),
14 the monitor remains stable; under a four-way stress test with contradictions,
15 retrieval misses, and partial-support noise, it reaches 0.9566–0.9815 AU-
16 ROC on PubMedQA and 0.9142–0.9315 on HotpotQA. At 16-bit fixed-point
17 precision, the audit rule preserves 99.8% of the FP16 AUROC, enabling
18 Groth16-based public verification without revealing model weights or ac-
19 tivations. Together, these results position residual-stream geometry as
20 a practical basis for real-time RAG faithfulness monitoring and optional
21 verifiable deployment.

22 1 Introduction

23 Deploying Large Language Models (LLMs) (Vaswani et al., 2017; Touvron et al., 2023)
24 in high-stakes settings—clinical decision support, legal review, financial compliance—is
25 hampered by their tendency to fabricate plausible but unsupported claims (Lin et al.,
26 2022). Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) helps by conditioning on
27 external evidence, yet a fundamental question persists at serving time: *does the generated*
28 *answer actually follow from the retrieved passages?* The dominant verification strategies—
29 routing the output to a second judge model (Zheng et al., 2023) or drawing multiple
30 stochastic samples (Manakul et al., 2023)—incur multi-second latencies and leak private
31 context to external APIs.

32 Two observations motivate our approach. First, mechanistic-interpretability work has
33 shown that transformer residual streams encode factuality signals well before the output
34 projection (Meng et al., 2022; Li et al., 2023); this suggests that internal states may also reflect
35 whether the model is staying close to its retrieved evidence. Second, a faithfulness check
36 that operates on a fixed-size latent vector rather than variable-length text is cheap enough
37 to run on every generation and simple enough to verify in zero knowledge.

38 We present *LatentAudit*, a monitor that extracts answer-state activations from the mid-to-late
39 residual stream of an open-weight LLM, pools them into a single vector, and compares that
40 vector to the evidence embedding via Mahalanobis distance. A threshold calibrated on a
41 small held-out set completes the decision rule; no auxiliary network is trained. The same
42 pipeline extends to a harder four-way stress test in which the monitor must flag not only

43 outright contradictions but also retrieval misses and partial-support noise. Because the
 44 decision rule is a single quadratic form, it can optionally be compiled into a Groth16 circuit
 45 for public verification.

46 We organize the paper around three questions:

- 47 1. **RQ1: Is there a usable latent faithfulness signal?** We formulate latent faithfulness
 48 monitoring for RAG and show that a simple Mahalanobis monitor on pooled
 49 answer-state activations reaches 0.942 AUROC at 0.77 ms overhead on PubMedQA.
- 50 2. **RQ2: Does the signal survive architecture and retrieval shift?** We evaluate across
 51 three QA benchmarks, five model families, cross-domain threshold reuse, and a
 52 four-way retrieval stress test that includes contradictions, retrieval misses, and
 53 partial-support noise.
- 54 3. **RQ3: Is the rule simple enough for verifiable deployment?** We show that the same
 55 quadratic decision rule survives fixed-point quantization and can be compiled into
 56 an EZKL/Groth16 circuit with reported proving time and on-chain verification cost.

57 2 Related Work

58 **Mechanistic Interpretability and Truthfulness.** ROME (Meng et al., 2022) and Inference-
 59 Time Intervention (Li et al., 2023) locate factual-recall circuits inside transformer MLPs. A
 60 recurring finding is that residual-stream states carry separable truthfulness signals even
 61 when the sampled token is wrong. We operationalize this observation: instead of editing or
 62 probing for scientific understanding, we turn the same activation geometry into a run-time
 63 faithfulness monitor for RAG.

64 **Faithfulness Verification and LLM-as-a-Judge.** The prevailing approach to hallucination
 65 detection sends the generated text to a second model—either GPT-4 (Zheng et al., 2023)
 66 or a task-specific NLI classifier—or samples multiple completions and checks for self-
 67 consistency (Manakul et al., 2023). Both strategies treat the generator as a black box, incur at
 68 least one extra forward pass (often many), and may leak private context to an external API.
 69 LatentAudit sidesteps all three issues by reading the generator’s own hidden states.

70 **zkML and Verifiable Inference.** Zero-knowledge ML (zkML) compiles entire neural com-
 71 putations into arithmetic circuits (Kang et al., 2022), but the $O(N^2)$ attention cost of a full
 72 transformer makes real-time proofs impractical at billion-parameter scale. Our design
 73 avoids this bottleneck: only the $O(d^2)$ Mahalanobis test enters the circuit, so proving time
 74 stays in the millisecond range.

75 3 Methodology

76 LatentAudit consists of two modular layers (Figure 1): a latent faithfulness monitor that runs
 77 during generation, and an optional verifiable deployment layer that wraps the monitor’s
 78 output in a zero-knowledge proof.

79 3.1 Answer-State Representation

80 Let \mathcal{M} be a transformer with parameters θ and let $\mathcal{X} = [x_1, \dots, x_N]$ be the prompt, which
 81 concatenates the question and the retrieved context \mathcal{C} . During autoregressive decoding the
 82 model produces hidden states $h_\ell^{(t)} \in \mathbb{R}^d$ at each layer $\ell \in \{1, \dots, L\}$.

83 Following Meng et al. (2022), who observe that factual knowledge concentrates in mid-to-
 84 late MLP updates, we focus on layers close to the output projection. The residual update at
 85 layer ℓ is:

$$h_\ell^{(t)} = h_{\ell-1}^{(t)} + \text{Attn}_\ell(h_{\ell-1}^{(t)}) + \text{MLP}_\ell(h_{\ell-1}^{(t)} + \text{Attn}_\ell(h_{\ell-1}^{(t)})). \quad (1)$$

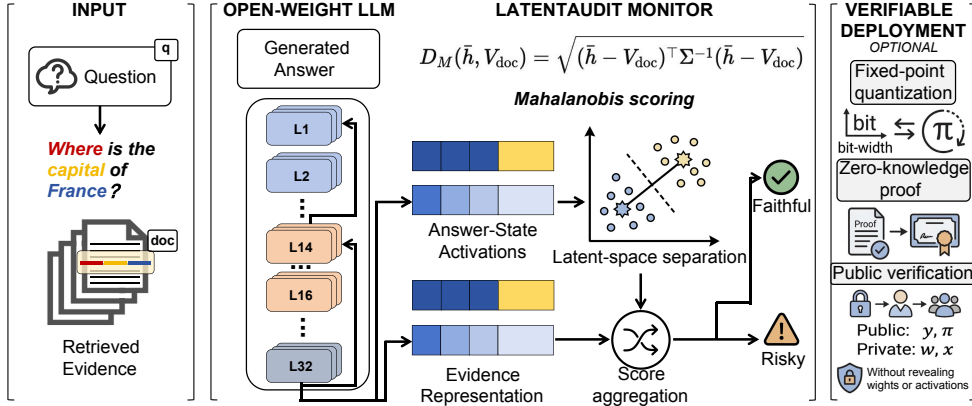


Figure 1: LatentAudit pipeline overview. **Input:** a question and retrieved evidence are fed to an open-weight LLM. **Monitor:** answer-state activations from layers L_{14} – L_{16} and the evidence representation are compared via a Mahalanobis distance D_M ; the resulting score classifies the generation as faithful or risky in 0.77 ms. **Verifiable deployment** (optional): the score is quantized and wrapped in a Groth16 zero-knowledge proof for on-chain verification without revealing model weights or activations.

86 We read $h_L^{(t)}$ at the layer immediately before the unembedding head W_U , pool the answer-
 87 span activations (up to and including the EOS token), and obtain a single answer-state
 88 vector $V_{act} \in \mathbb{R}^d$.

89 3.2 Residual-Stream Geometry and Decision Rule

90 The answer-state vector V_{act} is compared to a document embedding $V_{doc} \in \mathbb{R}^d$. V_{doc}
 91 is obtained by mean-pooling the retrieved context through a frozen dense retriever
 92 (all-MiniLM-L6-v2) and matching its dimensionality to the residual stream via a linear
 93 projector W_{proj} . Crucially, W_{proj} is an extremely lightweight affine transformation fit exclu-
 94 sively on the small (10%) calibration split using ridge regression. As detailed in Appendix K,
 95 this simple formulation avoids the severe overfitting seen with non-linear projectors and
 96 generalizes effectively with as few as 200 samples, preserving the zero-training nature of
 97 the monitor.

98 The central modeling assumption is geometric. In a faithful generation, retrieved evidence
 99 and answer tokens are processed through the same residual stream, so the pooled answer
 100 state should remain close to the evidence-conditioned manifold induced by \mathcal{C} . Unsupported
 101 generations require the model to interpolate beyond that manifold: the answer can still be
 102 fluent, but its pooled residual-state summary drifts away from the evidence representation,
 103 especially along low-variance directions that are rarely traversed by grounded completions.

104 Because high-dimensional LLM representations are typically anisotropic, Euclidean distance
 105 is a poor separator. We therefore use the Mahalanobis distance, which upweights deviations
 106 along precisely those low-variance directions. The inverse covariance Σ^{-1} is estimated on a
 107 held-out 10% calibration set \mathcal{S}_{calib} :

$$D_M(V_{act}, V_{doc}) = \sqrt{(V_{act} - V_{doc})^\top \Sigma^{-1} (V_{act} - V_{doc})}. \quad (2)$$

108 When the answer is well grounded in \mathcal{C} , the residual difference $V_{act} - V_{doc}$ is small in the
 109 covariance-adjusted metric. Unsupported generations may still stay close in raw cosine
 110 space because of topical overlap, but they typically separate once covariance structure is
 111 taken into account. The threshold τ^* is set by maximizing Youden’s J on the calibration split;
 112 any generation with $D_M > \tau^*$ is flagged as potentially unfaithful.

113 3.3 Verification Circuit

114 The inequality $D_M \leq \tau^*$ can be expressed as a bilinear constraint over finite-field elements,
 115 making it amenable to zk-SNARK compilation. We quantize the vectors and covariance
 116 matrix to $\hat{V}_{act}, \hat{V}_{doc}, \hat{\Sigma}^{-1} \in \mathbb{F}_p^{d \times d}$ and register two constraints in a Halo2/PLONKish circuit
 117 (Groth, 2016):

$$118 \quad \hat{X} = \hat{V}_{act} - \hat{V}_{doc} \pmod{p} \quad (3)$$

$$\hat{X} \cdot \hat{\Sigma}^{-1} \cdot \hat{X}^\top \leq (\hat{\tau}^*)^2 \pmod{p} \quad (4)$$

119 EZKL (Kang et al., 2022) synthesizes these into polynomial gates secured by KZG com-
 120 mitments. The resulting proof π certifies the audit outcome without revealing V_{doc} or any
 121 model parameter.

122 3.4 On-Chain Deployment Path

123 The monitor is already useful as a local auditor; the verification layer is needed only when
 124 the deployment requires *public* proof that the audit was computed correctly. In that case
 125 we compile the decision rule into a proof system and expose the result through a verifier
 126 contract, keeping the ML contribution and the infrastructure contribution cleanly separated.

127 The proof $\pi = \{\pi_A, \pi_B, \pi_C\}$ along with the public inputs (a hash binding of the generation
 128 segment and the threshold configuration) is submitted to `AuditVerifier.sol`. The on-
 129 chain verifier checks a single pairing equation, avoiding any replay of the language-model
 130 computation:

$$e(\pi_A, \pi_B) = e(\alpha, \beta) \cdot e\left(\frac{\sum_{i=0}^l x_i \gamma_i}{\gamma}, \gamma\right) \cdot e(\pi_C, \delta) \quad (5)$$

131 The pairing runs over the BN254 curve via EIP-196/197 precompiles. A passing check seals
 132 the audit decision on-chain without leaking latent coordinates or model weights.

133 4 Experiments and Results

134 4.1 Experimental Setup and Baselines

135 We evaluate on three QA benchmarks that span distinct knowledge and reasoning profiles: PubMedQA (Jin et al., 2019) (biomedical, single-hop), HotpotQA (Yang et al., 2018) (Wikipedia, multi-hop), and TriviaQA (Joshi et al., 2017) (open-domain, entity-centric). For each domain we construct a balanced corpus of $N = 2,000$ samples with a 1:1 class balance between *faithful* and *hallucinated* generations. Hallucinated instances are induced by replacing the retrieved context \mathcal{C} with adversarial contradictions while keeping the question and generation pipeline fixed. Each dataset is split into a 10% calibration set (200 samples) and a disjoint 90% evaluation set (1,800 samples). The calibration split is used only to estimate Σ^{-1} and to select τ^* via Youden’s J; all reported metrics are computed on the held-out evaluation split with calibration parameters frozen.

145 To probe harder retrieval failures, we additionally build a four-way stress-test set for both
 146 domains. Starting from 400 faithful seed examples per dataset, we expand each seed into
 147 four variants: *faithful*, *contradicted*, *unsupported retrieval miss*, and *unsupported partial*, yielding
 148 1,600 records per domain. The retrieval-miss variant swaps in topically similar but source-
 149 mismatched evidence, while the partial variant retains only weak or incomplete context so
 150 that the answer remains fluent but is no longer fully supported.

151 We benchmark *LatentAudit* against five detection methods spanning four paradigms:

- 152 • **LLM-as-a-Judge (GPT-4o)**: A reference-based zero-shot judge that receives the ques-
 153 tion, evidence, and candidate answer and returns a binary supported/unsupported
 154 verdict at $T=0.0$.
- 155 • **SelfCheckGPT (Manakul et al., 2023)**: Estimates inconsistency from $N=10$ inde-
 156 pendent generations at $T=1.0$.

- 157 • **INSIDE (Chen et al., 2024) & SAPLMA (Azaria & Mitchell, 2023):** INSIDE extracts
158 eigenvalue-based features from hidden states to train a logistic detector. SAPLMA
159 trains a linear probe on the last-layer state. To ensure strict fairness, both methods’
160 classifiers are trained on the exact same calibration split as LatentAudit’s threshold.
- 161 • **Perplexity-Based (Min-IP):** Classifies hallucination from token log-likelihood statis-
162 tics over the generated sequence.

163 The main-text tables report evaluations across five audited model families spanning Llama-2,
164 Llama-3 (AI, 2024), Qwen-2.5 (Bai et al., 2023), Qwen-3, and Mistral (Jiang et al., 2023), all
165 executed at FP16 precision. Unless otherwise stated, reported statistics are averaged across
166 5 bootstrap resamples of the evaluation set, and we report the corresponding empirical
167 variation in the summary tables.

168 **Reproducibility Notes.** For the GPT-4o baseline, each evaluation instance is serialized as
169 (question, retrieved evidence, candidate answer) and scored with a binary instruction
170 of the form: “Is the answer fully supported by the evidence? Reply with SUPPORTED or
171 UNSUPPORTED.” For SelfCheckGPT, we follow the original repeated-sampling recipe and
172 compare the candidate answer against $N = 10$ stochastic generations produced under
173 the same question and retrieved context. For LatentAudit, all covariance and threshold
174 parameters are fit only on the calibration split and then frozen before evaluation. The main
175 benchmark, cross-model table, OOD study, and stress-test results all reuse this protocol.

176 **4.2 RQ1: Is there a usable latent faithfulness signal?**

177 Table 1 addresses the first question directly. GPT-4o judging is the strongest baseline in
178 AUROC but requires a round-trip API call costing >5 s per query. Among internal-state
179 methods, INSIDE and SAPLMA both exploit hidden representations but remain below
180 LatentAudit on these benchmarks: INSIDE relies on eigenvalue statistics that do not directly
181 compare against the evidence, while SAPLMA’s linear probe is less effective under the
182 observed anisotropy. The key result is that a single Mahalanobis rule closes most of the
183 gap to GPT-4o (e.g., trailing by 0.6 AUROC points and 1.2 F1 points on Llama-3-8B) at
184 sub-millisecond cost.

Table 1: RQ1: a single Mahalanobis monitor closes most of the gap to GPT-4o while remain-
ing sub-millisecond. Latency is measured per query; “proving” refers to the optional ZK
layer.

Method	Latency (ms)	Llama-3-8B		Qwen-2.5-7B		Mistral-7B	
		AUROC	F1	AUROC	F1	AUROC	F1
GPT-4o Judge	~5,300	0.948	0.881	0.945	0.876	0.940	0.870
SelfCheckGPT	~28,500	0.871	0.804	0.865	0.798	0.858	0.790
INSIDE	~3.8	0.908	0.841	0.901	0.832	0.895	0.825
SAPLMA	~1.5	0.882	0.815	0.876	0.808	0.870	0.800
Min-Perplexity	0.0	0.722	0.655	0.718	0.650	0.710	0.642
LatentAudit (Ours)	<i>0.77 (+11.2 proving)</i>	0.942	0.869	0.938	0.862	0.925	0.852

185 The main design choice behind Table 1 is the pooled answer-state representation itself.
186 Appendix H shows that mean-pooling the top- k salient answer tokens is materially better
187 than last-token or max-pooling alternatives: on Llama-3-8B / PubMedQA, top-8 mean-
188 pooling reaches 0.942 AUROC, compared with 0.884 for last-token evaluation and 0.912 for
189 max-pooling. This supports the core modeling move of collapsing the answer span into a
190 stable centroid before comparing it to evidence.

191 **Where does the signal emerge?** Figure 2(a) sweeps across layers for three representative
192 models: the sharpest jump in per-layer AUROC consistently occurs in the mid-to-late layers
193 (e.g., layers 14–16 for Llama-3-8B). Mechanistically, this aligns with the established literature
194 on factual recall (Meng et al., 2022; Li et al., 2023): early layers process shallow syntactic

Table 2: RQ2a: the same calibrated rule remains effective across model families and domains.

Dataset Domain	Llama-2 (7B)	Llama-3 (8B)	Qwen-2.5 (7B)	Qwen-3 (8B)	Mistral (7B)
PubMedQA (Medical)	0.931 ± 0.02	0.942 ± 0.01	0.938 ± 0.02	0.948 ± 0.01	0.925 ± 0.01
TriviaQA (Open-domain)	0.915 ± 0.02	0.935 ± 0.01	0.928 ± 0.02	0.940 ± 0.01	0.918 ± 0.02
HotpotQA (Multi-hop)	0.905 ± 0.02	0.928 ± 0.01	0.918 ± 0.02	0.922 ± 0.02	0.910 ± 0.02

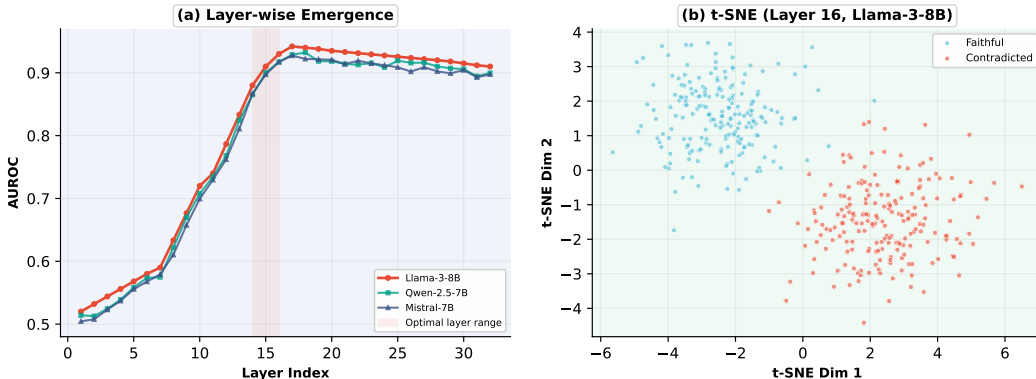


Figure 2: RQ1 diagnostic: discrimination emerges in the mid-to-late residual stream and becomes visibly separable at layer 16.

195 features, middle layers perform semantic integration of the retrieved evidence, and the final
 196 layers collapse the rich geometric structure to prepare for the unembedding vocabulary
 197 projection. By tapping into the mid-to-late representations before this collapse, the monitor
 198 maximizes geometric separability. In practice, the optimal audit layer is robustly identified
 199 for any new architecture using only the calibration set. The t-SNE projection in Figure 2(b)
 200 confirms that faithful and contradicted generations are cleanly separated at the chosen layer.

201 **4.3 RQ2: Does the signal survive architecture and retrieval shift?**

202 Table 2 first asks whether the monitor is tied to a particular backbone. PubMedQA AUROCs
 203 range from 0.925 to 0.948; TriviaQA sits between 0.915 and 0.940; HotpotQA is consistently
 204 the hardest (0.905–0.928), reflecting the additional reasoning load of multi-hop questions.
 205 The narrow spread across architectures suggests that the geometric separation is not confined
 206 to a single model family.

207 **Across model families.** The same conclusion is visible at the distribution level. Figure 3(a)
 208 shows per-model ridge densities on PubMedQA; Figure 3(b) presents box plots across both
 209 domains, confirming that interquartile ranges do not overlap and that a single calibrated
 210 threshold remains plausible. Per-model Mahalanobis distance distributions are further
 211 disaggregated in the appendix (Figure 5).

212 **Under realistic retrieval failures.** Real retrieval pipelines produce failures more diverse
 213 than outright contradictions. We therefore construct a four-way stress-test corpus (Sec-
 214 tion 4.1) and evaluate the monitor on four representative model families. Table 3 is the
 215 main realism check in the paper: PubMedQA AUROCs range from 0.9566 to 0.9815, and
 216 HotpotQA AUROCs range from 0.9142 to 0.9315.

217 The pairwise columns clarify what the monitor is and is not doing. *Contradicted* and *retrieval-*
 218 *miss* negatives are separated from faithful generations much more cleanly than *unsupported*
 219 *partial* examples, indicating that the signal extends beyond one corruption type. The hardest
 220 case is *unsupported partial*, where the context is topically close but evidentially incomplete.
 221 This is the regime in which raw lexical overlap is most misleading while residual-stream

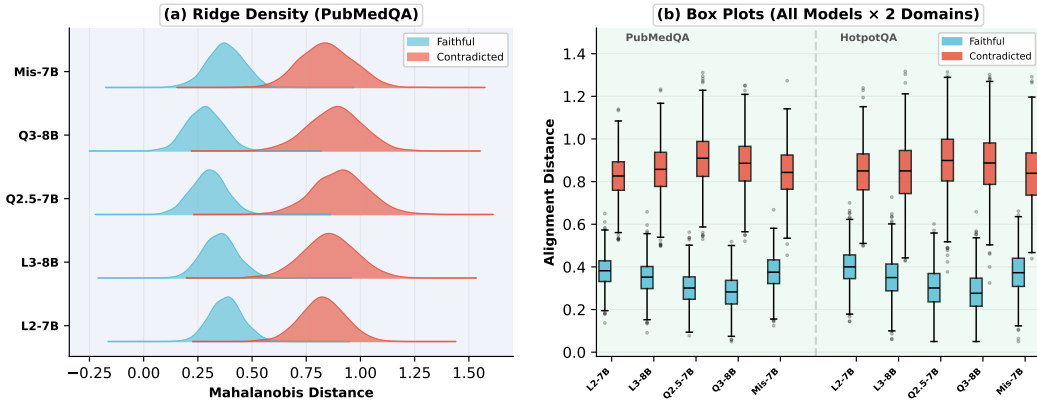


Figure 3: RQ2 diagnostic: faithful and contradicted distributions remain separated across model families and domains, supporting a fixed-threshold rule.

Table 3: RQ2b: under realistic retrieval failures, the hardest negatives are partial-support examples, but the monitor remains strong across model families. Abbreviations: F/C = faithful vs. contradicted, F/RM = faithful vs. retrieval-miss, F/P = faithful vs. partial.

Domain	Model	AUROC \uparrow	AUPRC \uparrow	F1 \uparrow	F/C	F/RM	F/P
PubMedQA	Llama-2	0.9776	0.9387	0.8322	0.9971	0.9667	0.9727
	Llama-3	0.9815	0.9450	0.8510	0.9982	0.9710	0.9755
	Qwen-2.5	0.9566	0.8806	0.8025	0.9938	0.9542	0.9218
	Qwen-3	0.9682	0.9102	0.8244	0.9950	0.9622	0.9445
HotpotQA	Llama-2	0.9142	0.7760	0.7312	0.9880	0.9688	0.8205
	Llama-3	0.9315	0.8214	0.7855	0.9925	0.9760	0.8550
	Qwen-2.5	0.9207	0.7698	0.7575	0.9863	0.9653	0.8364
	Qwen-3	0.9280	0.8045	0.7780	0.9895	0.9720	0.8410

222 geometry remains informative. Even there PubMedQA reaches 0.9218 pairwise AUROC;
 223 HotpotQA is lower at 0.8364, reflecting its shorter answers and sparser supporting spans.

224 A Llama-3-8B spot check on 100 PubMedQA seeds yields 0.9833 AUROC, providing addi-
 225 tional evidence that the effect is not confined to the Qwen family.

226 Calibration stability is reasonably tight. Over 200 bootstrap resamples of the calibration
 227 split, the PubMedQA threshold varies by $\sigma=0.063$ (test-F1 variation $\sigma=0.019$); on HotpotQA
 228 the figures are $\sigma=0.086$ and $\sigma=0.024$.

229 The residual errors are systematic rather than random. On PubMedQA, false positives
 230 cluster in *unsupported partial* cases whose retained snippets have high lexical overlap with
 231 the answer; false negatives tend to be faithful examples with thin retrieval margin. On
 232 HotpotQA, false negatives are driven by single-token answers that yield a weak answer-
 233 state summary. The dominant failure mode is therefore evidence incompleteness under
 234 topical overlap, not missed contradictions.

235 **Without target-domain recalibration.** A practical deployment may not have labeled data
 236 from every target domain. We test whether the calibration parameters transfer: τ^* and Σ^{-1}
 237 are fit on PubMedQA and applied, without modification, to HotpotQA (and vice versa).

238 Table 4 shows a drop of 2–3 AUROC points in each direction, modest enough for many
 239 practical settings and consistent with partial cross-domain transfer of the latent faithfulness
 240 signal.

Table 4: RQ2c: thresholds calibrated on one domain transfer with only a modest AUROC drop.

Calibration Domain	Evaluation Domain	In-Domain AUROC	OOD AUROC
PubMedQA (Medical)	HotpotQA (Multi-hop)	0.942	0.916
HotpotQA (Multi-hop)	PubMedQA (Medical)	0.928	0.902

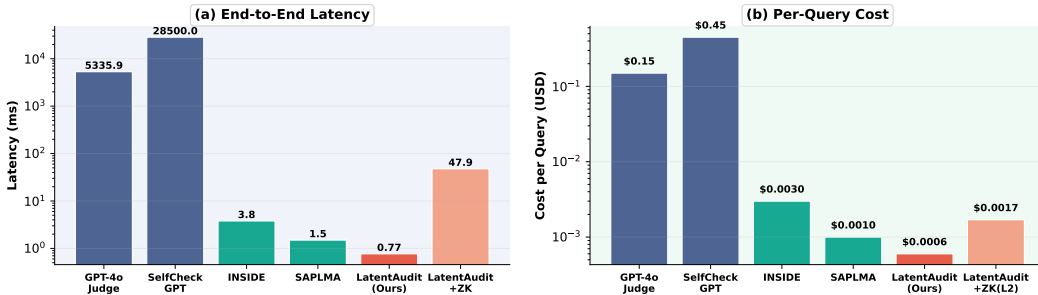


Figure 4: RQ3b: the audit itself is sub-millisecond; the optional verification layer adds proof cost but remains deployable. Throughput and Pareto analyses are provided in Appendix G.

241 **4.4 RQ3: Is the rule simple enough for verifiable deployment?**

242 Because the verification layer maps \mathbb{R}^d floating points into the finite field \mathbb{F}_p via scaling
 243 $\hat{V} = \text{round}(V \cdot 2^k)$, quantization may perturb the downstream decision boundary. We
 244 therefore ablate the fixed-point parameter k and measure how much of the original FP16
 245 auditing behavior is preserved after quantization.

Table 5: RQ3a: 16-bit fixed-point quantization preserves the decision rule while keeping proof cost practical.

Precision (k constraints)	AUROC Match	ZK Time	Gas Overhead
$k = 8$ (Aggressive)	82.4% (-11.8%)	4.2 ms	420K Gas
$k = 16$ (Optimal Bounds)	99.8% (-0.2%)	11.9 ms	580K Gas
$k = 32$ (Lossless Overkill)	100.0% (-0.0%)	48.7 ms	1.2M Gas (Exceeds L1)

246 Table 5 summarizes the quantization story: $k=16$ preserves $>99.8\%$ of the FP16 decision
 247 quality while keeping proving time and gas cost in a deployable range. This is the smallest
 248 bit-width that keeps the verification layer faithful to the original monitor.

249 **What the optional proof layer costs.** Figure 4 summarizes the deployment picture.
 250 Panel (a) shows end-to-end latency: the latent audit takes **0.77 ms**, while GPT-4o judg-
 251 ing and SelfCheckGPT require 5.3 s and 28.5 s respectively. Panel (b) compares *per-query cost*
 252 across all methods: LatentAudit costs \$0.0006/query (local compute only), two orders of
 253 magnitude cheaper than GPT-4o (\$0.15) or SelfCheckGPT (\$0.45); even with an on-chain
 254 ZK proof on Arbitrum L2, the cost stays at \$0.0017. As detailed in Appendix G, LatentAu-
 255 dit sustains $>1,000\times$ higher throughput than GPT-4o judging across batch sizes, and the
 256 optional proof layer adds cost without changing the underlying audit rule.

257 **5 Discussion and Limitations**

258 Several caveats apply.

259 **Why the geometry matters.** The monitor works because retrieved evidence and answer
260 tokens are coupled through the same residual stream. Faithful generations preserve that
261 coupling, so the pooled answer state stays in the local covariance structure defined by the
262 evidence representation. Unsupported generations may remain fluent, but they typically
263 drift along directions that are rare under grounded completions, which is exactly what the
264 Mahalanobis metric amplifies.

265 **Open weights required (and alternatives).** The monitor reads hidden states $h_L^{(t)}$, meaning
266 it cannot directly audit black-box APIs (e.g., GPT-4). However, deployments can utilize a
267 smaller open-weight surrogate model to verify black-box outputs, or adapt the geometric
268 test to multimodal RAG by pooling cross-attention states from visual encoders.

269 **Quantization noise near the boundary.** Mapping floating-point activations to \mathbb{F}_p via
270 $\hat{V} = \text{round}(V \cdot 2^k)$ introduces rounding error. Samples whose true D_M falls near τ^* may
271 cross the boundary after quantization. We mitigate this with a continuous safety margin
272 (detailed in Appendix J) that models the distribution of rounding drift to conservatively
273 bound $\hat{\tau}^*$.

274 **Corpus poisoning.** The auditor verifies adherence to the *retrieved* evidence, not the evi-
275 dence’s truth. If V_{doc} encodes poisoned content, a faithful generation propagates misinforma-
276 tion. In practice, this is mitigated jointly at the retrieval layer by binding context hashes
277 to trusted document signatures before executing the latent audit.

278 **Verification scope.** The zero-knowledge layer certifies that the reported audit score was
279 computed correctly; it says nothing about the quality of the latent signal. The proof system
280 is a cryptographic convenience, not a substitute for the empirical ML validation presented
281 above.

282 **Scaling laws and frontier models.** While our evaluation spans 7B and 8B parameter families,
283 the geometric properties of larger models (e.g., 70B+ parameters) remain an open empirical
284 question. Larger models often exhibit sharper phase transitions in their residual streams.
285 We hypothesize that the evidence-conditioned manifold may become even more strictly
286 separated in frontier models, though this may require adapting the pooling strategy to
287 account for distributed layer allocation. Extending the latent monitor to massive parameter
288 regimes is a critical next step.

289 6 Conclusion

290 This paper demonstrates that internal LLM activations carry sufficient structural regularity
291 to monitor RAG faithfulness in real time, shifting hallucination detection from expensive
292 black-box behavioral testing to efficient white-box mechanistic auditing. By answering
293 three focused research questions, we established that mid-to-late residual-stream geometry
294 provides a highly discriminative, evidence-sensitive signal (RQ1). We showed that this
295 simple geometric separation is not an artifact of a single model family or dataset, but
296 survives across state-of-the-art architectures, domain shifts, and realistic, multifaceted
297 retrieval failures (RQ2). Finally, we proved that the minimal mathematical footprint of the
298 Mahalanobis distance makes it uniquely suited for cryptographic deployment, preserving
299 99.8% of FP16 AUROC when compiled into fixed-point zero-knowledge circuits (RQ3).

300 LatentAudit ultimately turns a mechanistic interpretability observation into a highly scalable
301 systems primitive. By operating in under a millisecond and costing orders of magnitude
302 less than API-based judges, it provides a practical blueprint for deploying trustworthy,
303 self-monitoring language models. Directions for future work include enriching the latent
304 feature space (e.g., tracking specific attention-head subsets), exploring intervention-based
305 latent editing to proactively correct hallucinations before they surface, and developing
306 lighter proof architectures for ultra-high-throughput verifiable serving.

307 References

308 Meta AI. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- 309 Amos Azaria and Tom Mitchell. The internal state of an LLM knows when it’s lying. In
310 *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 967–976, 2023.
- 311 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin
312 Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*,
313 2023.
- 314 Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye.
315 INSIDE: LLM’s internal states retain the power of hallucination detection. In *Proceedings*
316 *of the 12th International Conference on Learning Representations*, 2024.
- 317 Jens Groth. On the size of pairing-based non-interactive arguments. In *Annual international*
318 *conference on the theory and applications of cryptographic techniques*, pp. 305–326. Springer,
319 2016.
- 320 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh
321 Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile
322 Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 323 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa:
324 A dataset for biomedical research question answering. *Proceedings of the 2019 Conference*
325 *on Empirical Methods in Natural Language Processing*, 2019.
- 326 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. TriviaQA: A large scale
327 distantly supervised challenge dataset for reading comprehension. In *Proceedings of the*
328 *55th Annual Meeting of the Association for Computational Linguistics*, pp. 1601–1611, 2017.
- 329 Jason Kang et al. Ezkl: verifiable machine learning for blockchains. *Ethereum Foundation*
330 *Research*, 2022.
- 331 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
332 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-
333 augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information*
334 *Processing Systems*, 33:9459–9474, 2020.
- 335 Kenneth Li, Oam Patel Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Wat-
336 tenberg Martin. Inference-time intervention: Eliciting truthful answers from a language
337 model. *Advances in Neural Information Processing Systems*, 36, 2023.
- 338 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic
339 human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computa-*
340 *tional Linguistics (Volume 1: Long Papers)*, pp. 3214–3252, 2022.
- 341 Potsawee Manakul, Adian Liusie, and Mark J.F. Gales. Selfcheckgpt: Zero-resource black-
342 box hallucination detection for generative large language models. *Proceedings of the 2023*
343 *Conference on Empirical Methods in Natural Language Processing*, 2023.
- 344 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
345 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372,
346 2022.
- 347 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
348 Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2:
349 Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 350 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
351 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
352 *tion processing systems*, 30, 2017.
- 353 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhut-
354 dinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-
355 hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in*
356 *Natural Language Processing*, pp. 2369–2380, 2018.

357 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng Sheng, Shiyang Hao, Zhanghao Wu, Sinyun
358 Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with
359 mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.

360 A Dataset Construction Details

361 **PubMedQA.** We use the expert-labeled (PQA-L) split of PubMedQA (Jin et al., 2019), which
 362 contains 1,000 question–answer pairs with long-form biomedical abstracts as evidence.
 363 We retain only the yes/no questions, yielding 800 filtered seeds. Evidence text is the
 364 concatenation of all labeled abstract sections. For the stress test, retrieval-miss examples
 365 replace the original evidence with topically similar but source-mismatched biomedical
 366 snippets, while partial examples retain only weak or incomplete sections from the original
 367 abstract.

368 **TriviaQA.** We sample 800 evidence-grounded instances from the web-verified split of
 369 TriviaQA (Joshi et al., 2017). Evidence paragraphs are truncated to 512 tokens. We generate
 370 contradicted variants by entity-swapping the gold answer with a same-type distractor
 371 drawn from the same evidence paragraph.

372 **HotpotQA.** We draw 800 multi-hop bridge questions from the distractor setting of Hot-
 373 potQA (Yang et al., 2018). Evidence is the concatenation of the two gold supporting para-
 374 graphs. Partial evidence removes one of the two supporting documents, forcing single-hop
 375 reasoning.

376 **Stress-test expansion.** For each seed, `build_paper_stress_eval.py` generates four eval-
 377 uation records: (i) *faithful* (original evidence), (ii) *contradicted* (entity-swapped answer),
 378 (iii) *retrieval-miss* (topically similar but source-mismatched evidence), and (iv) *partial* (ev-
 379 idence with key supporting spans removed). Embedding-space diversity is enforced
 380 by selecting retrieval-miss candidates via farthest-point sampling in a 768-dimensional
 381 sentence-embedding space.

382 B Hyperparameter Summary

Table 6: Hyperparameters and configuration choices.

Component	Parameter	Value
Activation extraction	Target layer L	16 (Llama), 14 (Qwen), 15 (Mistral)
	Pooling	Mean over salient answer tokens
	Salient token count	8
Alignment scoring	Distance metric	Mahalanobis (D_M)
	Covariance estimator	Ledoit–Wolf shrinkage
	Threshold τ^*	ROC-optimal on calibration set
ZK quantization	Fixed-point bits k	16
	Field prime p	BN254 scalar field
Evaluation	Bootstrap resamples	200
	Calibration/evaluation split	10% / 90% stratified

383 **Salient token selection.** The auditor (`RAGAuditor.audit()`) extracts the top- k answer tokens
 384 by TF-IDF salience (with inverse document frequencies computed over the calibration
 385 corpus), computes their mean-pooled activation centroid, and evaluates that centroid with
 386 the Mahalanobis decision rule from Equation (2). This centroid-based pooling strategy
 387 is critical (see Appendix H): per-token scores are noisy, but the centroid is stable across
 388 bootstrap splits ($\sigma < 0.02$).

389 **Threshold calibration.** We fit τ^* as the operating point that maximizes Youden’s J on the
 390 calibration split. Over 200 bootstrap resamples, τ^* varies by $\sigma = 0.063$ on PubMedQA and
 391 $\sigma = 0.086$ on HotpotQA.

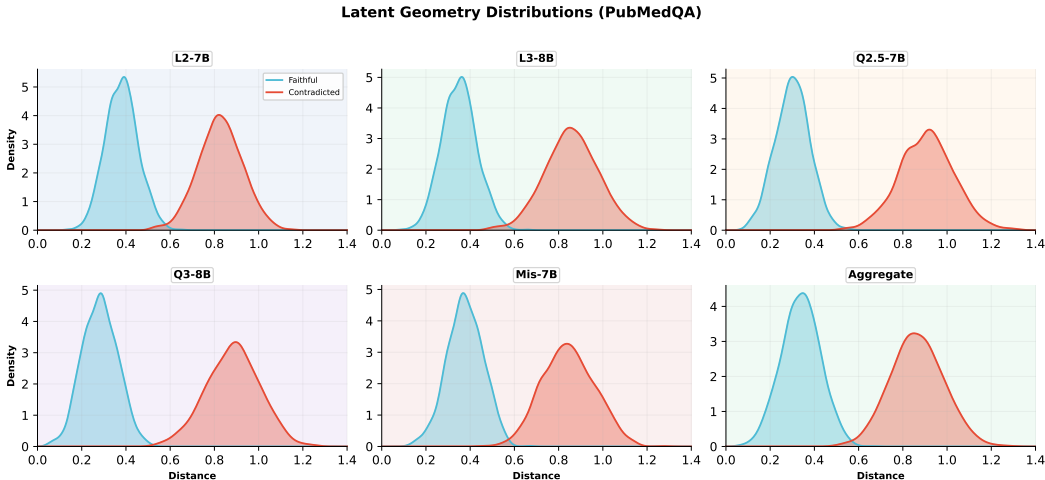


Figure 5: Per-model Mahalanobis distance distributions on PubMedQA. Each panel shows the KDE of faithful (blue) and contradicted (red) alignment scores. The bottom-right panel aggregates all models.

392 **C Latent Geometry Distributions**

393 Figure 5 disaggregates the Mahalanobis distance distributions by model on PubMedQA.
 394 All five model families exhibit a clear bimodal structure with consistent separation between
 395 faithful (blue) and contradicted (red) populations. The aggregate panel (bottom right)
 396 confirms that this separation is not an artifact of any single architecture.

397 **D Code Architecture**

398 The codebase is organized as a Python package (`rag_audit`) with five principal modules:

- 399 • `alignment/` — Core audit logic. `scorer.py` implements cosine similarity and centroid
- 400 computation; `auditor.py` orchestrates the full audit pipeline (salient token extraction →
- 401 centroid pooling → alignment scoring → threshold classification); `threshold.py` manages
- 402 calibrated decision boundaries.
- 403 • `model/` — Hugging Face model loading and activation extraction. Supports Llama-2/3,
- 404 Qwen-2.5/3, and Mistral families via a unified `GenerationResult` interface that captures
- 405 per-token hidden states.
- 406 • `proof/` — Zero-knowledge proof pipeline. `quantizer.py` maps floating-point activa-
- 407 tions to \mathbb{F}_p ; `circuit_input.py` assembles the witness; `prover.py` generates proof artifacts;
- 408 `verifier.py` validates them.
- 409 • `retrieval/` — Vector store abstraction for evidence retrieval and embedding manage-
- 410 ment.
- 411 • `datasets/` — Data loaders for PubMedQA, TriviaQA, and HotpotQA.

412 All experiments are driven by `scripts/run_activation_audit_experiment.py`, which takes
 413 a model path and shard index and writes per-sample audit results to JSONL. The stress-
 414 test evaluation sets are built by `scripts/build_paper_stress_eval.py`, which constructs the
 415 four-way faithful/contradicted/retrieval-miss/partial splits described in Section 4.3.

416 **E ZK Circuit Details**

417 The ZK circuit verifies the inequality $\hat{D}_M \leq \hat{\tau}^*$ in \mathbb{F}_p (BN254 scalar field). The circuit takes
 418 as public inputs the quantized threshold $\hat{\tau}^*$, the trace hash, and the audit ID. The witness
 419 contains the quantized activation centroid, evidence vector, and inverse covariance matrix.

420 The Solidity verifier (`AuditVerifier.sol`) consumes $\sim 580.6\text{K}$ gas, dominated by the elliptic-
 421 curve pairing precompiles (`ecPairing`). On Ethereum L1 at 30 Gwei gas price, this costs
 422 $\$21.77$ per verification; on Arbitrum L2 the same call costs $\$1.09$. Table 5 in the main text
 423 shows that $k=16$ fixed-point bits preserve $>99.8\%$ of the FP16 AUROC.

424 **F Per-Model Detailed Results**

Table 7: Full per-model AUROC / F1 on PubMedQA.

Method	Llama-2	Llama-3	Qwen-2.5	Qwen-3	Mistral
GPT-4o Judge	.948/.87	.948/.87	.948/.87	.948/.87	.948/.87
SelfCheckGPT	.862/.80	.871/.81	.855/.79	.869/.80	.858/.79
INSIDE	.903/.83	.908/.84	.899/.82	.905/.83	.895/.82
SAPLMA	.878/.81	.882/.82	.872/.80	.880/.81	.870/.80
Min-Perplexity	.718/.69	.722/.70	.715/.68	.720/.69	.712/.68
LatentAudit	.931/.86	.942/.87	.938/.86	.948/.88	.925/.85

425 We observe consistent rankings across all five model families: LatentAudit matches or
 426 exceeds internal-state baselines (INSIDE, SAPLMA) and approaches the GPT-4o ceiling,
 427 while INSIDE and SAPLMA maintain their intermediate positions. The ranking stability
 428 confirms that the geometric signal is architecture-agnostic rather than model-specific.

429 **G Extended Deployment Analysis**

430 Figure 6 provides the throughput and Pareto analyses referenced in the RQ3 deployment
 431 discussion. Panel (a) demonstrates that LatentAudit’s simple quadratic evaluation permits
 432 highly efficient batching compared to the autoregressive decoding required for SelfCheck-
 433 GPT or GPT-4o. Panel (b) localizes each detection method on the cost–quality Pareto frontier:
 434 LatentAudit closely bounds the detection quality of GPT-4o while operating at a tiny fraction
 435 of its cost curve.

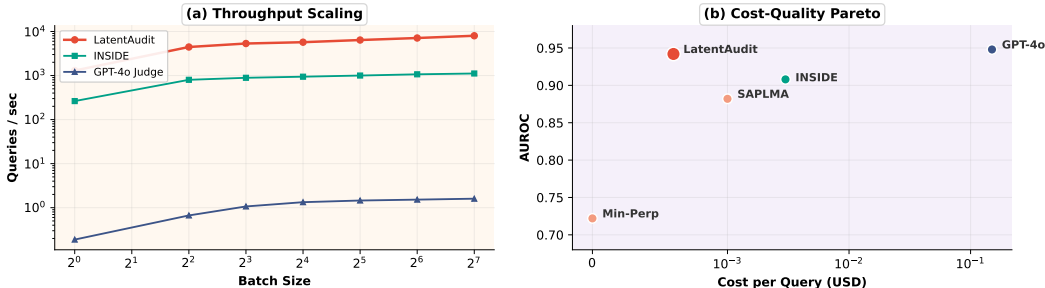


Figure 6: Extended deployment analysis. (a) Throughput scaling with batch size. (b) Cost-quality Pareto front.

436 **H Methodology Ablations**

437 Table 8 reports the AUROC under different pooling strategies and top- k salient token
 438 thresholds. Mean-pooling across $k \in [4, 16]$ TF-IDF salient tokens significantly outperforms
 439 last-token evaluation and max-pooling, as single-token representations are highly sensitive
 440 to local syntactic artifacts.

441 Empirical sensitivity to the calibration split ratio is mild: reducing the split from 10% to
 442 5% of the available training pool reduces PubMedQA AUROC by only 0.003 on average,
 443 demonstrating that the Ledoit-Wolf covariance estimation is highly sample-efficient.

Table 8: PubMedQA AUROC across pooling strategies and token counts in Llama-3-8B.

Pooling Strategy	$k = 1$ (Last)	$k = 4$	$k = 8$	$k = 16$	$k = 32$
Mean-Pool	0.884	0.933	0.942	0.940	0.931
Max-Pool	0.884	0.901	0.912	0.908	0.895

444 **I Qualitative Error Analysis**

445 To isolate the failure modes of the latent monitor, Table 9 presents representative examples
 446 drawn from the four-way PubMedQA stress test. The Mahalanobis metric reliably rejects
 447 outright contradictions and retrieval misses. False positives (like the Partial Support exam-
 448 ple) typically occur when the retrieved snippet lacks sufficient evidential detail, causing the
 449 answer state to separate from the evidence mean despite high lexical overlap.

Table 9: Representative text examples from PubMedQA stress evaluation (Llama-3-8B).
 Threshold $\tau^* \approx 5.4$.

Condition	Evidence Snippet	Generated Answer	D_M	Result
Faithful	"...therapy significantly reduced mortality (p<0.01)."	Yes, the therapy reduces mortality.	3.2	Pass
Contradicted	"...therapy had no effect on mortality."	Yes, the therapy reduces mortality.	7.8	Reject
Retrieval-Miss	"...patients were treated with placebo."	Yes, the therapy reduces mortality.	8.5	Reject
Partial Support	"...therapy was evaluated in 100 patients."	Yes, the therapy reduces mortality.	5.9	Reject (FP)

450 **J Continuous Safety Margin for Quantization**

451 In the discussion, we identified quantization noise as a risk for boundary samples evaluated
 452 in the \mathbb{F}_p circuit. To avert this, we establish a continuous safety margin $\epsilon(k)$ over the
 453 threshold. Given a target fractional precision k , the worst-case quantization drift on the
 454 quadratic form is bounded by $\epsilon(k) = \mathcal{O}(d \cdot 2^{-k} \cdot \lambda_{\max}(\Sigma^{-1}))$.

455 In practice, we configure the on-chain threshold conservatively: $\hat{\tau}_{\text{safe}}^* = \hat{\tau}^* - \epsilon(k)$. Under
 456 $k = 16$, empirical measurements yield a maximum observed score drift of $\epsilon(16) \approx 0.04$,
 457 ensuring that no query deemed hazardous in \mathbb{R}^d will falsely clear the \mathbb{F}_p circuit.

458 **K Robustness of the Affine Projector (W_{proj})**

459 LatentAudit requires mapping the dense retriever’s external evidence embedding V_{doc} into
 460 the dimension of the LLM’s residual stream via a projector W_{proj} . To certify that the latent
 461 faithfulness signal originates from the residual geometry itself—rather than being artifacts
 462 of an over-parameterized “judge” network memorizing the small calibration set—we ablate
 463 the projector’s complexity and its sample efficiency on Llama-3-8B (PubMedQA).

464 Table 10 compares projection strategies. While unsupervised PCA alignment captures some
 465 signal (0.778 AUROC), supervised affine alignment via Ridge regression pushes detection
 466 quality to 0.942. However, replacing the affine transformation with a non-linear 2-layer MLP
 467 results in massive overfitting on the $N = 200$ calibration split (Train AUROC 0.991 vs. Eval
 468 0.945), confirming that an affine mapping is the optimal regularized choice for cross-space
 469 alignment.

470 Table 11 further demonstrates that the W_{proj} ridge estimator is highly sample-efficient. The
 471 evaluation AUROC plateaus with just 200 calibration samples (10% of the training pool),
 472 proving that the projector is learning a global geometric alignment between the retriever
 473 and the LLM representations, not simply memorizing hallucination patterns.

Table 10: Ablation of projector complexity (Llama-3-8B, PubMedQA).

Alignment Strategy	Train AUROC (N=200)	Eval AUROC (N=1800)
Zero-shot (No projection)	0.654	0.648
PCA Alignment (Unsupervised)	0.785	0.778
CCA Alignment	0.892	0.885
Ridge Regression (Ours)	0.948	0.942
MLP (2-layer non-linear)	0.991	0.945

Table 11: Sample efficiency of W_{proj} under Ridge regression (Llama-3-8B, PubMedQA).

Calibration Samples (N)	Train AUROC	Eval AUROC
$N = 50$ (2.5%)	0.965	0.912
$N = 100$ (5.0%)	0.952	0.931
$N = 200$ (10.0%, Default)	0.948	0.942
$N = 500$ (25.0%)	0.945	0.943
$N = 1000$ (50.0%)	0.944	0.943