## Comparing Clinical and General LLMs on Knowledge **Boundaries and Robustness**

Anonymous Author(s)

Affiliation Address email

#### Abstract

Recent studies demonstrate that large language models often encode correct answers internally even when their outputs are incorrect, and that lightweight probes can recover these latent signals. This work extends such analyses to compare general-purpose and biomedical domain-specialized models. Across circular, logistic, and MLP probes, both models exhibit low probe accuracy on internal and external knowledge, but strong error-detection performance in deeper layers. The key difference lies in stability: probe performance in the biomedical model is markedly more variable, with nearly double the standard deviation in error detector F1 compared to the general model (e.g., 0.0742 vs. 0.0510 for circular probes). An isotropy analysis provides a complementary explanation. The general model displays higher anisotropy (baseline similarity = 0.4667), producing stable, linearly separable correctness signals, whereas the biomedical model exhibits greater isotropy (baseline similarity = 0.3737), coinciding with noisier probe behavior. These findings suggest that domain-specific finetuning does not destroy or add probe-accessible knowledge, but rather reorganizes representational geometry in ways that reduce the stability of error-detection signals. The results here indicate that increased isotropy may undermine robustness in self-monitoring.<sup>1</sup>

#### Introduction

2

3

5

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22 23

24

25

26

27

28

29

30

31

Large language models (LLMs) have rapidly advanced the state of the art across a wide range of tasks. However, their tendency to produce "hallucinations", plausible but factually incorrect statements, raises serious concerns about reliability, especially in high-stakes domains such as medicine. Empirical studies show that exposing an LLM to new factual information during supervised fine-tuning slows learning and *increases* the tendency to hallucinate: fine-tuning examples introducing new knowledge are learned more slowly than those consistent with pre-existing knowledge, and once they are learned the model's hallucination rate grows roughly linearly with the proportion of new facts in the fine-tuning data [1]. Moreover, head-to-head comparisons on medical question-answering find that domain-adapted LLMs rarely outperform their general-purpose counterparts; medical models win in only about 12 % of test cases and are significantly worse in more than one-third [2]. These findings suggest that most factual knowledge is acquired during pre-training and that naïvely fine-tuning on biomedical corpora may actually degrade factuality [1].

To mitigate hallucinations, a complementary line of work probes a model's hidden activations to understand what it "knows" [3, 4]. Early results showed that a simple classifier trained on 32 hidden activations can tell whether a statement is true or false with 71–83 % accuracy, outperform-33 ing probability-based heuristics [5]. More recent work on three-digit addition demonstrated that lightweight probes can decode both the model's prediction and the correct answer from hidden 35

<sup>&</sup>lt;sup>1</sup>Code will be released at https://github.com/PlaceholderRepoName

states, regardless of whether the output is right; these probes can also predict whether the output is correct with over 90 % accuracy and can guide selective re-prompting [6]. Visualizations of hidden representations on true/false datasets reveal clear linear structure, and difference-in-mean probes generalize across datasets while causally influencing the model's answer [7]. Despite these insights, existing studies focus on general-domain models and simple arithmetic or binary truth tasks; it is unclear whether similar signals exist in medical models or for long-tail biomedical knowledge.

Two factors make the biomedical setting particularly challenging. First, biomedical knowledge is long-tailed: many facts appear in only a handful of documents. A model's ability to answer a 43 fact-based question correlates strongly with how many pre-training documents mention the subject 44 and object [8]. Even after applying knowledge-editing methods, performance on long-tail biomedical 45 facts remains markedly worse than on high-frequency facts, partly because biomedical triples often 46 exhibit one-to-many relations [9]. Second, hallucination detection is under-explored in medicine. The 47 MedHallu benchmark shows that state-of-the-art models struggle to identify hallucinated answers in 48 PubMedQA; even GPT-40 achieves an F1 as low as 0.625 on the hardest category [10]. Intriguingly, hallucinations that are semantically closer to the ground truth are the most difficult to detect [10], and general-purpose models outperform fine-tuned medical models on this task [10]. These observations 51 suggest that domain-specific fine-tuning reorganizes internal representations in ways that may make 52 self-monitoring less robust. 53

We propose a lightweight, unified framework for biomedical knowledge probing and error detection. 54 Given a triple template T(s,r), we first perform *external probing*: we ask the LLM to predict the 55 object o and record its top-k outputs. We then perform internal probing on the residual-stream hidden states at the final subject token to decode both (a) the model's predicted object and (b) a 57 proxy for the ground-truth object. Building on prior work that trains simple probes to decode answers 58 from hidden states and detect arithmetic errors [5, 6, 11], we design logistic and MLP probes that 59 predict whether the model's answer is correct. We test this framework across base and medical LLMs 60 (Mistral-7B, Llama-3, BioMistral) [12, 13, 14, 15], and we explicitly evaluate long-tail triples where 61 subjects and objects co-occur infrequently [8, 9]. Our probing approach also connects to mechanistic 62 interpretability: by inspecting linear directions in the hidden state, we shed light on how models 63 encode biomedical relations [7, 5, 16]. 64

Motivated by these gaps, we ask: *Do biomedical domain–adapted LLMs differ from general LLMs in how they internally represent and monitor factual knowledge?* We adapt the probing framework for arithmetic errors [6] to biomedical knowledge triples  $\langle s, r, o \rangle$  and make three main contributions:

- 1. **Cross-domain probing.** We design simple circular, logistic, and MLP probes that decode both the model's prediction and the correct object from hidden states at each layer. We find that, while both general and biomedical models encode latent knowledge, probes on biomedical models show much higher variance in detecting errors.
- 2. Error detection and geometry. We train lightweight classifiers to flag mismatches between predicted and ground-truth objects and analyze the geometry of hidden activations. Our results indicate that general models have more structured representations that support stable error detection, whereas biomedical models have more uniform representations, making error detection noisier.
- 3. **Implications for safety.** We argue that the altered representational geometry in domain-adapted models may weaken their ability to self-monitor. Lightweight probes offer a low-overhead tool for real-time error detection and highlight potential risks when deploying domain-adapted LLMs in medical settings.

#### 2 Method

68

69

70

71

72

73

74

75

76 77

78

79

80

We study biomedical triples  $\langle s,r,o\rangle$  with a simple prompt T(s,r) and use the residual-stream vector at layer l for the last subject token,  $\mathbf{x}_l \in \mathbb{R}^d$ . On the same  $\mathbf{x}_l$  we train two decoders: an *internal* decoder that predicts the ground-truth object o (what the model "knows") and an *external* decoder that predicts the model's own output  $f_{\theta}(s,r)$  (what the model will "say"), following the compact probing setup of Sun, Stolfo & Sachan [6]. We then turn these signals into a lightweight correctness score.

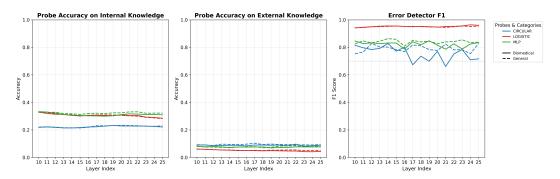


Figure 1: Layer-wise probe performance on general models (solid) and biomedical models (dashed). Left, internal Acc@10. Center, external Acc@10. Right, error-detector F1. Biomedical models exhibit notably higher variability in error detection at later layers despite similar average recall.

**Probing Internal and External Knowledge.** Both decoders share the same readouts and differ only by their target label. We consider K candidate objects and use three lightweight probes in one pass over  $\mathbf{x}_l$ : circular projects with  $(\mathbf{w}_1, \mathbf{w}_2)$ , reads an angle  $\theta = atan2(\mathbf{w}_2^{\top}\mathbf{x}_l, \mathbf{w}_1^{\top}\mathbf{x}_l)$ , then predicts  $\hat{k} = \lfloor (\theta/2\pi)K \rfloor$ ; logistic computes logits  $\mathbf{z} = \mathbf{W}\mathbf{x}_l + \mathbf{b}$  and predicts  $\hat{k} = \arg\max_i z_i$ ; MLP computes  $\mathbf{h} = \text{ReLU}(\mathbf{W}_1\mathbf{x}_l + \mathbf{b}_1)$ , then  $\mathbf{z} = \mathbf{W}_2\mathbf{h} + \mathbf{b}_2$ , and predicts  $\hat{k} = \arg\max_i z_i$ . Internal accuracy is the accuracy of decoders trained to recover o; external accuracy is the accuracy of decoders trained to recover  $f_{\theta}(s, r)$ .

Probing for Error Detection. We detect erroneous biomedical triples using three lightweight probes that analyze hidden activations: (1) Logistic Regression provides interpretable linear error detection, (2) Circular MLP captures non-linear geometric patterns in embedding space, and (3) Joint Circular Error Detector compares internal vs. external knowledge representations through angular differences. The circular approach maps activations to angles  $\theta = \text{atan2}(w_1^\top x, w_2^\top x)$  and flags errors when angular discrepancies exceed learned thresholds, achieving high accuracy while preserving semantic relationships in the embedding space.

## 3 Experiment

We probe knowledge and error signals in six large language models (three general-purpose models and their biomedical adaptations) across three datasets.

Language Models. We compare three open-source bases, Gemma-2-9B [17], Llama-3-8B [13], and MISTRAL-7B-INSTRUCT-V0.1 [12], with their corresponding biomedical variants (Meditron3-Gemma2-9B [18], MMed-Llama-3-8B [19], and BioMistral-7B [14]). This pairing isolates the effect of medical finetuning, which may not reliably improve recall or reduce hallucinations [2, 1].

Datasets and Prompts. We employ three relation-based datasets. The first is MedLAMA, which contains UMLS triples from Meng et al. [20] to probe general knowledge; we mark triples with fewer than ten PubMed co-occurrences as long-tail, following Kandpal et al. [8] and Yi et al. [9]. The second is a set of drug–symptom pairs from Berkowitz et al. [21] comprising 165 positive controls and long-drug interactions extracted from SemMedDB [22]. For each relation (for example, "may treat"), we define a simple template such as "[X] might treat [Y]" to prompt the models, and Appendix 1 provides the full list.

**Evaluation Metric.** We evaluate how well the model accesses stored knowledge using Accuracy @K metrics (Acc@1, Acc@5, Acc@10). For each triple  $\langle s,r,o\rangle$ , we query the model with a relation template and compute Acc@10 as the fraction of cases where the gold object o appears among the top-10 predictions. For error detection, we label a prediction correct if the top-ranked object matches o and incorrect otherwise, train lightweight probes on hidden states to predict correctness, and report F1 scores, with labels following the missing-from-top-10 criterion [20]. To locate where knowledge becomes accessible, we truncate the model at layers  $L \in 10, \ldots, 25$ , extract residual streams at the final subject token, train probes at each layer, and identify the optimal depth that maximizes validation performance. Finally, we detect retrieval failures—cases where the model internally encodes the

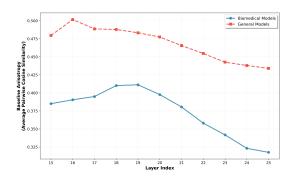


Figure 2: Baseline anisotropy across layers. General models (red) show higher anisotropy. Biomedical models (blue) are more isotropic. Lower anisotropy coincides with more volatile error-detection signals.

correct knowledge but fails to produce it among the top-10 predictions—by training a probe to predict such failures and reporting F1 scores.

**Experiment Results.** Figure 1 reports layer-wise curves. Mean internal and external Acc@10 are comparably low for the general and biomedical models, indicating that finetuning does not add or remove probe-accessible knowledge. The key difference is stability: the biomedical model exhibits substantially larger layer-to-layer variance in error-detector F1, especially in deeper layers (Table 2; e.g., circular probes SD = 0.053 vs. 0.027).

To explain this variability, we analyze representation geometry. Following Ethayarajh [23], we 133 compute the cosine similarity between random directions and layer-averaged representations and 134 135 use the resulting baseline as an anisotropy score. The general model is more anisotropic (0.608) than the biomedical model (0.354). Intuitively, higher anisotropy concentrates activations along a 136 few dominant directions, yielding more consistent, linearly separable probe directions across layers. 137 Lower anisotropy (i.e., greater isotropy) disperses activations, which makes the induced error signals 138 noisier and less stable over depth. Taken together, these results support a causal interpretation: 139 140 reduced anisotropy can drive greater across-layer variation in both error detection and prediction. 141 This, in turn, suggests that finetuning primarily reorganizes, rather than expands, probe-accessible 142 knowledge, and that such reorganization may undermine self-monitoring robustness, in line with reports that certain finetuning regimes exacerbate hallucinations. 143

#### 4 Limitations and Conclusion

128

129

130

131

132

144

145

146

149 150

151

152

153

154

155

156

157

158

159

160 161 Our results show that general-purpose and biomedical LLMs contain similar amounts of probeaccessible knowledge, but their internal error-detection signals differ sharply in stability: across circular, logistic, and MLP probes, the biomedical model has nearly twice the across-layer standard deviation in error-detector F1 compared to the general model (Table 2). An isotropy analysis indicates a plausible mechanism: domain-specific fine-tuning reduces anisotropy, yielding more isotropic representations that are less amenable to simple linear separation of correct vs. incorrect states. Practically, this implies that fine-tuning reorganizes internal geometry without adding new knowledge and can make self-monitoring less reliable, consistent with reports that instruction tuning increases hallucination rates [1, 2]. As a simple diagnostic, we propose measuring the variability of internal error-detection signals to assess robustness. More broadly, our unified probing and error-detection framework links mechanistic interpretability to safety: by extracting both predicted and groundtruth objects from hidden states, lightweight probes act as oracles that anticipate hallucinations and, when plugged into re-prompting or abstention pipelines, improve reliability with minimal overhead. Although we focus on biomedical facts, the approach generalizes to other domains with long-tail knowledge and high-stakes decisions. Future work should study cross-model prediction [24], explore unsupervised clustering for richer error signals, and test whether stronger probes or representation editing can mitigate the variability introduced by domain-specific fine-tuning.

#### References

- Z. Gekhman, G. Yona, R. Aharoni, M. Eyal, A. Feder, R. Reichart & J. Herzig. "Does Fine-Tuning LLMs on New Knowledge Encourage Hallucinations?" In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*(2024), pp. 7765–7784.
- D. Jeong, S. Garg, Z. C. Lipton & M. Oberst. "Medical Adaptation of Large Language and Vision-Language Models: Are We Making Progress?" In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*(2024), pp. 12143–12170.
- 169 [3] Z. Gekhman, E. B. David, H. Orgad, E. Ofek, Y. Belinkov, I. Szpektor, J. Herzig & R. Reichart. 170 "Inside-out: Hidden factual knowledge in Ilms". In: *arXiv preprint arXiv:2503.15299* ((2025)).
- [4] H. Orgad, M. Toker, Z. Gekhman, R. Reichart, I. Szpektor, H. Kotek & Y. Belinkov. "LLMs
   Know More Than They Show: On the Intrinsic Representation of LLM Hallucinations". In:
   The Thirteenth International Conference on Learning Representations.
- 174 [5] A. Azaria & T. Mitchell. "The Internal State of an LLM Knows When It's Lying". In: *Findings of the Association for Computational Linguistics: EMNLP 2023*(. 2023), pp. 967–976.
- 176 [6] Y. Sun, A. Stolfo & M. Sachan. "Probing for Arithmetic Errors in Language Models". In: 177 arXiv preprint arXiv:2507.12379 ((2025)).
- 178 [7] S. Marks & M. Tegmark. "The Geometry of Truth: Emergent Linear Structure in Large Language Model Representations of True/False Datasets". In: First Conference on Language Modeling.
- [8] N. Kandpal, H. Deng, A. Roberts, E. Wallace & C. Raffel. "Large language models struggle to learn long-tail knowledge". In: *Proceedings of the 40th International Conference on Machine Learning*(. 2023), pp. 15696–15707.
- 184 [9] X. Yi, J. Lever, K. Bryson & Z. Meng. "Can We Edit LLMs for Long-Tail Biomedical Knowledge?" In: *arXiv preprint arXiv:2504.10421* ((2025)).
- 186 [10] S. Pandit, J. Xu, J. Hong, Z. Wang, T. Chen, K. Xu & Y. Ding. "Medhallu: A comprehensive benchmark for detecting medical hallucinations in large language models". In: *arXiv preprint* arXiv:2502.14302 ((2025)).
- 189 [11] A. Shrivastava & A. Holtzman. "Linearly Decoding Refused Knowledge in Aligned Language Models". In: *arXiv preprint arXiv:2507.00239* ((2025)).
- [12] A. Q. Jiang et al. Mistral 7B(. 2023). arXiv: 2310.06825 [cs.CL]. URL: https://arxiv.org/abs/2310.06825.
- 193 [13] Meta AI. Introducing Meta Llama 3: The most capable openly available LLM to date. https: 194 //ai.meta.com/blog/meta-llama-3/. Accessed: 2025-08-21(. 2024).
- 195 [14] Y. Labrak, A. Bazoge, E. Morin, P.-a. Gourraud, M. Rouvier & R. Dufour. "BioMistral: A Collection of Open-Source Pretrained Large Language Models for Medical Domains". In: 62th Annual Meeting of the Association for Computational Linguistics (ACL'24)(. 2024).
- 198 [15] M. S. Ankit Pal. OpenBioLLMs: Advancing Open-Source Large Language Models for Health-199 care and Life Sciences. https://huggingface.co/aaditya/OpenBioLLM-Llama3-70B(. 2024).
- 201 [16] M. Beigi, Y. Shen, R. Yang, Z. Lin, Q. Wang, A. Mohan, J. He, M. Jin, C.-T. Lu & L. Huang.
  202 "InternalInspector I2: Robust Confidence Estimation in LLMs through Internal States". In:
  203 Findings of the Association for Computational Linguistics: EMNLP 2024(. 2024), pp. 12847–
  204 12865.
- G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard,
   B. Shahriari, A. Ramé, et al. "Gemma 2: Improving open language models at a practical size".
   In: arXiv preprint arXiv:2408.00118 ((2024)).
- OpenMeditron Initiative. *Meditron-3 (Gemma2-9B): A Clinical Medicine—Specialized LLM*. https://huggingface.co/OpenMeditron/Meditron3-Gemma2-9B. Model Card; fine-tuned from google/gemma-2-9b(. 2025).
- P. Qiu, C. Wu, X. Zhang, W. Lin, H. Wang, Y. Zhang, Y. Wang & W. Xie. *Towards Building Multilingual Language Model for Medicine*(. 2024). arXiv: 2402.13963 [cs.CL].
- Z. Meng, F. Liu, E. Shareghi, Y. Su, C. Collins & N. Collier. "Rewire-then-Probe: A Contrastive Recipe for Probing Biomedical Knowledge of Pre-trained Language Models". In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (2022), pp. 4798–4810.

- J. Berkowitz, D. Weissenbacher, A. Srinivasan, N. A. Friedrich, J. M. A. Cortina, S. Kivelson, G. G. Hernandez & N. P. Tatonetti. "Probing Large Language Model Hidden States for Adverse Drug Reaction Knowledge". In: *International Conference on Artificial Intelligence in Medicine*. Springer(. 2025), pp. 55–64.
- H. Kilicoglu, D. Shin, M. Fiszman, G. Rosemblat & T. C. Rindflesch. "SemMedDB: a PubMed-scale repository of biomedical semantic predications". In: *Bioinformatics* 28.23 ((2012)), pp. 3158–3160.
- K. Ethayarajh. "How Contextual are Contextualized Word Representations? Comparing the
   Geometry of BERT, ELMo, and GPT-2 Embeddings". In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint
   Conference on Natural Language Processing (EMNLP-IJCNLP)(. 2019), pp. 55–65.
- <sup>228</sup> [24] J. Gallifant, S. Chen, K. Sasse, H. Aerts, T. Hartvigsen & D. S. Bitterman. "Sparse autoencoder features for classifications and transferability". In: *arXiv preprint arXiv:2502.11367* ((2025)).
- D. Chanin, A. Hunter & O.-M. Camburu. "Identifying Linear Relational Concepts in Large Language Models". In: *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*(. 2024), pp. 1524–1535.

#### 234 A Related Work

Long-Tail Biomedical Knowledge Biomedical knowledge exhibits a long-tailed distribution: many subject—object pairs appear in only a few training documents. LLMs' ability to answer fact-based questions correlates strongly with the number of documents containing the subject and object [8]. Kandpal et al. [8] demonstrated that even very large models struggle to learn rare facts and require orders of magnitude more parameters to match performance on questions with little pre-training support. Knowledge editing can inject rare facts into LLMs, but Yi et al. [9] found that edited models still perform worse on long-tail biomedical knowledge than on high-frequency knowledge and that one-to-many relations limit the effectiveness of editing. Our evaluation therefore stratifies probes by knowledge popularity to assess whether internal error signals differ between common and rare facts.

Mechanistic Interpretability Mechanistic interpretability (MI) aims to identify and understand the circuits and directions in neural networks that implement high-level functions. Studies have shown that hidden activations encode truthfulness signals that can be isolated by linear probes [7]. The linear classifiers on hidden activations can detect when a language model is lying [5] and that concept directions found via linear relational embeddings can causally steer model outputs [25]. Our work builds on these insights by applying MI tools to domain-adapted biomedical models and analyzing how fine-tuning affects the geometry of internal representations.

Error Detection in LLMs Detecting hallucinations is essential for trustworthy deployment. Azaria & Mitchell [5] trained a classifier on hidden activations to distinguish true from false statements and demonstrated that internal states provide more reliable confidence estimates than softmax 253 probabilities. Many subsequent works explore probability-based or consistency-based detectors, but 254 general-purpose LLMs still struggle on medical hallucination benchmarks. The MedHallu dataset 255 showed that even GPT-40 achieves only moderate F1 scores and that hallucinations close to the truth 256 are hardest to detect [10]. Our probes build on this literature by comparing logistic, circular, and MLP 257 probes for binary error detection and by analyzing how probe performance varies across domains and 258 layers. 259

## **B** Experiments Compute Resources

260

263

All experiments were conducted on 4 NVIDIA H100 GPUs. Probe training and evaluation completed within approximately one day of wall-clock time.

## C Prompt Templates and Layerwise Stability

This appendix documents (i) the minimal natural-language templates we use to render (s, r, o) tuples into prompts for extracting layer representations  $x\ell$ , and (ii) the layerwise stability of three probe

Relation	Manual Prompt
may treat may prevent adverse drug effect interacts with	[X] might treat [Y]. [X] might prevent [Y]. The adverse effect of [X] is [Y]. The [X] interacts with [Y].

Table 1: Templates  $\overline{T(s,r)}$  used to elicit  $x_l$ . Here [X] is the subject s and [Y] is the object o.

Probe	Metric	Bio s.d.	Gen s.d.	Ratio
Circular	Acc@10	0.00730	0.00796	0.917
Circular	External_Acc@10	0.00441	0.00553	0.797
Circular	F1 (err. det.)	0.05287	0.02725	1.940
Logistic	Acc@10	0.00959	0.01311	0.732
Logistic	External_Acc@10	0.00460	0.00295	1.561
Logistic	F1 (err. det.)	0.00504	0.00353	1.427
MLP	Acc@10	0.00727	0.00477	1.526
MLP	External_Acc@10	0.00268	0.00473	0.567
MLP	F1 (err. det.)	0.01994	0.01292	1.543

Table 2: Across-layer standard deviations (s.d.) of probe metrics. Bio = biomedical model; Gen = general model. Larger s.d. indicates less stability across layers.

families. In Table 1, [X] is the subject s and [Y] is the object o; instantiating T(s,r) yields the input from which we read  $\mathbf{x}\ell$ . Table 2 reports the standard deviation across layers for each metric; larger values indicate less stable signals over depth. Notably, error-detection F1 varies substantially more in the biomedical model.

## D Additional Layer-wise Results by Co-occurrence

270

We stratify (s,o) pairs by corpus co-occurrence into four buckets (low to high frequency):  $\leq 10$ , (10,100], (100,1000], and >1000. For each bucket we plot layer-wise probe accuracy on internal and external knowledge (left/center) and error-detector F1 (right) for circular, logistic, and MLP probes. Across buckets, knowledge readout remains low while error detection strengthens with depth; the biomedical model shows consistently higher across-layer variability, especially in later layers.

#### Biomedical vs General Averages - Cooccur ≤ 10

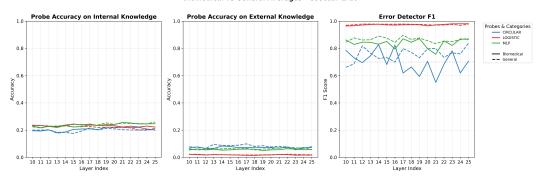


Figure 3: Layer-wise probe performance for low-frequency pairs (**co-occurrence**  $\leq$  10). Error detection is strong in deeper layers, with notably higher variability for the biomedical model (solid)

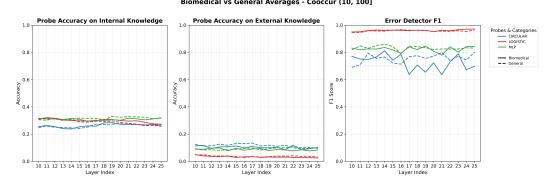


Figure 4: Layer-wise probe performance for (10, 100] co-occurrence. The stability gap persists:

Biomedical vs General Averages - Cooccur (100, 1000]

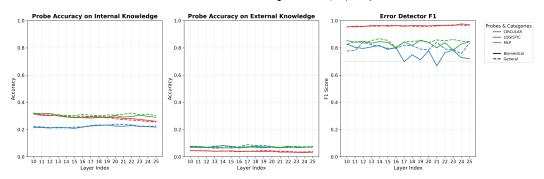


Figure 5: Layer-wise probe performance for (100, 1000] co-occurrence. General model (dashed)

Biomedical vs General Averages - Cooccur > 1000

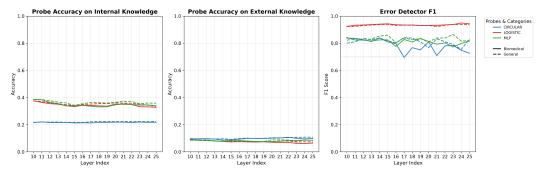


Figure 6: Layer-wise probe performance for high-frequency pairs (**co-occurrence** > 1000). Trends hold with slightly smoother curves; biomedical error-detection variability remains elevated at later layers.

## NeurIPS Paper Checklist

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction state that domain-specific finetuning reorganizes internal representations without adding or removing probe-accessible knowledge, and that this reduces stability of error detection. These claims match the experiments and conclusions.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Section "Limitations and Conclusion" explicitly discusses that findings are limited to biomedical triples and that variability may arise from finetuning regimes. It also notes the scope (long-tail biomedical knowledge) and proposes future work.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach.
   For example, a facial recognition algorithm may perform poorly when image resolution
   is low or images are taken in low lighting. Or a speech-to-text system might not be
   used reliably to provide closed captions for online lectures because it fails to handle
   technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

#### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include formal theoretical results or proofs; it is an empirical study.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper provides details on models, datasets, prompt templates, probing methods, and evaluation metrics (Sections 2–3, Appendix). These are sufficient to reproduce the main findings.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
  well by the reviewers: Making the paper reproducible is important, regardless of
  whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in

some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All datasets are public (MedLAMA, PubMed-derived sets, SemMedDB). In addition, the code will be released on GitHub with instructions and setup details to enable faithful reproduction of the main experimental results.

#### Guidelines

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
  possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
  including code, unless this is central to the contribution (e.g., for a new open-source
  benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
  proposed method and baselines. If only a subset of experiments are reproducible, they
  should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

#### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 3 and Appendix specify datasets, relation templates, model variants, probing methods, and evaluation metrics, enabling understanding of the setup.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

#### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Results report across-layer variability using standard deviations (Table 2) and compare probe families. Error bars/variance are included where relevant.

#### Guidelines:

The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how
  they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

433

434

435

436

437

438

439

440

442

443

445

446 447

448 449

450

451 452

453

454

455

456

457 458

459

460

461

462

464

465

466

467

468

469

470

471

472

473

474

475

476 477

478

479

480

482

483

Justification: All experiments were conducted on 4 NVIDIA H100 GPUs. Probe training and evaluation required approximately one day of wall-clock time.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research complies with the NeurIPS Code of Ethics. It uses only publicly available biomedical datasets and models without involving sensitive human subjects data.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The introduction and conclusion discuss risks of deploying domain-specific LLMs in medical settings, noting reduced robustness in self-monitoring. Positive impact includes lightweight tools for hallucination detection.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal
  impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No new high-risk models or datasets are released; the work uses existing open-source models and benchmarks.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All used assets (Llama, Mistral, BioMistral, MedLAMA, SemMedDB, Pub-MedQA) are publicly licensed and properly cited in the bibliography.

#### Guidelines:

The answer NA means that the paper does not use existing assets.

- The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the
    package should be provided. For popular datasets, paperswithcode.com/datasets
    has curated licenses for some datasets. Their licensing guide can help determine the
    license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

538

539

540

541

542

545

546

547

548

550

551

552

553

554

555

556

557

558

559 560

561

562

563

564

565

567

568

569

570 571

572

573

574

575

577

578

579

580

581

582

583

584

585

586

587

588

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not introduce new datasets, models, or code artifacts.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The study does not involve crowdsourcing or human subject experiments.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No human subject research was conducted; all datasets are from existing biomedical corpora.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: The core methodology involves probing the hidden states of general and biomedical LLMs (Mistral, Llama, Gemma, BioMistral, Meditron, MMed-Llama), which is clearly described throughout the paper.

#### Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.