

---

# InterpDetect: Interpretable Signals for Detecting Hallucinations in Retrieval-Augmented Generation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Retrieval-Augmented Generation (RAG) reduces hallucinations by incorporating  
2 external knowledge, yet models still produce content inconsistent with retrieved  
3 evidence. Effective detection requires separating external context from parametric  
4 knowledge, which prior methods conflate. We find hallucinations often emerge  
5 when later-layer FFNs inject excessive parametric knowledge into the residual  
6 stream. To address this, we explore *external context scores* and *parametric knowl-*  
7 *edge scores*, mechanistic features derived from Qwen3-0.6b across layers and  
8 attention heads. Using these signals, we train lightweight classifiers and evaluate  
9 against LLMs and detection baselines. Furthermore, classifiers trained on Qwen3-  
10 0.6b signals generalize to GPT-4.1-mini responses, demonstrating the potential  
11 of `proxy-model` evaluation. Our results show that mechanistic signals offer  
12 efficient, generalizable predictors for hallucination detection in RAG.<sup>1</sup>

## 13 1 Introduction

14 Large language models (LLMs) excel at tasks such as question answering and summarization [1, 2],  
15 yet they frequently generate *hallucinations*—outputs that are factually incorrect or unsupported [3].  
16 Retrieval-Augmented Generation (RAG) reduces hallucinations by grounding responses in external  
17 knowledge [4], but models still produce outputs inconsistent with retrieved content [5, 6, 7], making  
18 reliable detection essential for domains such as healthcare and finance. Recent advances in RAG hal-  
19 lucination detection follow two directions. Corpus-based approaches (e.g., RAGTruth, LettuceDetect,  
20 Mu-SHROOM) provide fine-grained supervision and cross-lingual benchmarks [8, 9, 10]. Mechanis-  
21 tic approaches (e.g., ReDeEP, LRP4RAG) analyze internal activations to disentangle external context  
22 from parametric knowledge, showing that overactive FFNs and mis-weighted attention often drive  
23 hallucinations [11, 12]. We follow the line of ReDeEP work [11] by introducing **External Context**  
24 **Score (ECS)** and **Parametric Knowledge Score (PKS)** as predictive features for hallucination  
25 detection. ECS measures grounding in retrieved content, while PKS captures FFN-driven parametric  
26 knowledge injection. Together, they decompose generation into external vs. internal contributions.  
27 Using Qwen3-0.6b, we compute ECS and PKS across layers and attention heads, then train regression-  
28 based classifiers for hallucination detection. We further demonstrate `proxy-model` evaluation:  
29 classifiers trained on Qwen3-0.6b transfer effectively to larger models (e.g., GPT-4-mini), offering  
30 scalable, low-cost detection.

31 Our contributions are threefold:

- 32 1. We extend the ReDeEP framework for ECS/PKS computation with a fully open-sourced imple-  
33 mentation built on TransformerLens [13], enabling any GPT-like model without modifying the  
34 underlying model library.

---

<sup>1</sup>Code and data: <https://anonymous.4open.science/r/interpretability-hallucination-experiments-75D3>.

- 35 2. We conduct a systematic evaluation of multiple regression-based classifiers, identifying the optimal  
 36 model for hallucination prediction.  
 37 3. We demonstrate that leveraging a 0.6b-parameter model as a proxy allows effective and economical  
 38 computation, facilitating practical application to large-scale, production-level models.

39 **2 Methodology**

40 Figure 1 outlines our methodology. We first run a standard RAG pipeline to generate responses.  
 41 The model’s *parametric knowledge* is treated as internal knowledge, elicited by prompting with  
 42 the query. For each context–response span, we compute ECS from attention heads and PKS from  
 43 feed-forward layers, capturing external vs. internal contributions. After validating their correlation  
 44 with hallucination labels, these scores serve as features for binary classifiers, whose span-level outputs  
 45 are aggregated into response-level predictions.

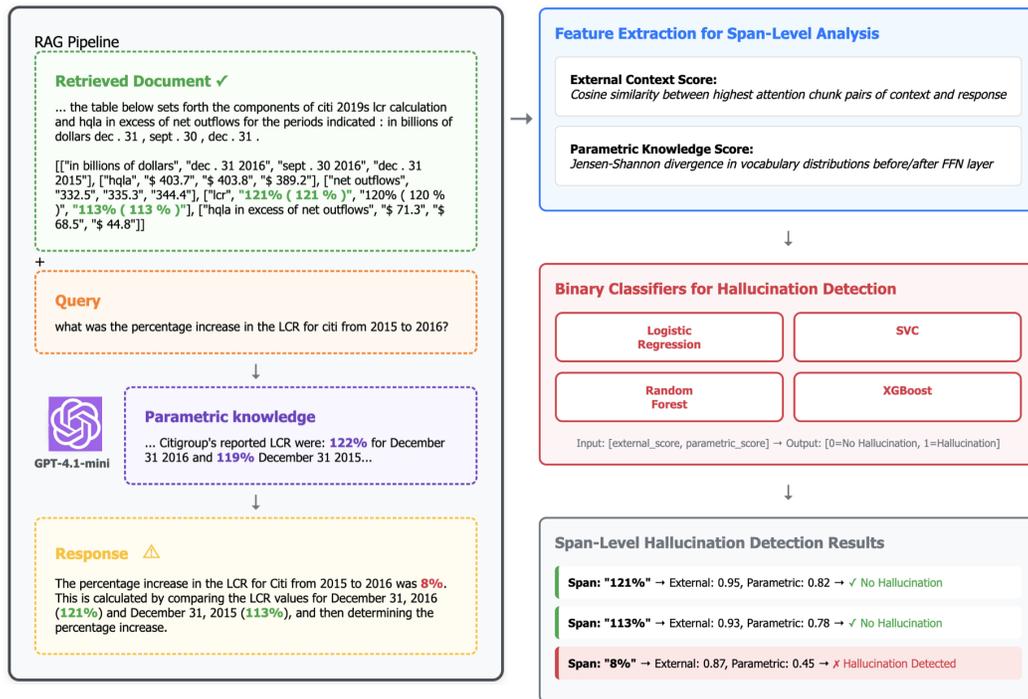


Figure 1: Pipeline for Span-Level Hallucination Detection

46 **Data Curation** We use the FinQA dataset from RagBench [14], which contains 12.5k training  
 47 and 2.29k test instances. Each instance pairs an S&P 500 financial report with a grounded question  
 48 and an LLM-generated response. For our study, we regenerate responses with Qwen3-0.6b (used  
 49 for interpretability scores) and GPT-4.1 mini to test cross-model generalization. Our data pipeline  
 50 consists of three steps (Appendix A shows an example):

- 51 • **Response Generation:** We subsample 3,000 training instances from the original dataset and  
 52 generate responses separately with Qwen3-0.6b and GPT-4.1 mini.  
 53 • **Labeling:** We label hallucinated spans using LettuceDetect [9], refined with two LLM judges  
 54 (llama-4-maverick-17b-128e-instruct and gpt-oss-120b for Qwen3-0.6b-generated response; claude-  
 55 sonnet-4 for GPT-4.1-mini-generated response to reduce family bias). Majority voting is applied at  
 56 the response level, and samples are retained if at least one judge confirms the LettuceDetect label,  
 57 yielding 1,852 instances for analysis.  
 58 • **Chunking:** To compute accurate ECS and PKS, we operate at the span level by chunking documents  
 59 (using documents\_sentences from FinQA) and splitting responses into individual sentences.

60 **Mechanistic Metrics** Our work draws on concepts from mechanistic interpretability research.  
 61 Specifically, we leverage TransformerLens, an open-source library that provides access to internal  
 62 model parameters—such as attention heads, feed-forward network layers (FFNs), and residual  
 63 streams—in GPT-like models. For a detailed description of the architecture, we refer the reader  
 64 to the original TransformerLens work [13]. We primarily use TransformerLens to compute two  
 65 metrics: the *External Context Score* and the *Parametric Knowledge Score*, following the definitions  
 66 in ReDeEP [11]. While the original ReDeEP framework supports both token-level and chunk-level  
 67 hallucination detection, computing scores at the token level is computationally expensive and does  
 68 not fully capture context. Therefore, we restrict our calculations at the chunk level. Details about the  
 69 computation is given in Appendix B.

70 **Classifier for Hallucination Detection** Hallucination detection is formulated as a binary classifica-  
 71 tion task. As input features, we use ECS computed for each attention head and layer, together with  
 72 PKS computed for each layer. Predictions are generated at the span level and can subsequently be  
 73 aggregated to obtain response-level results. Training details are provided in Appendix C.

### 74 3 Experiments

75 **Correlation Analysis** We examine mechanistic signals for their correlation with hallucinations  
 76 in RAG outputs. Comparing ECS values between truthful and hallucinated responses (Figure 2a)  
 77 shows that hallucinated responses rely less on external context. Pearson correlation analysis using  
 78 inverse hallucination labels confirms that higher ECS corresponds to lower hallucination likelihood  
 79 (Figure 2b). While copying-head proxies (`0V_copying_score`) offer interpretable, low-cost ap-  
 80 proximations [15], we did not observe strong correlations in Qwen3-0.6b, so we include all layers  
 81 and heads for ECS and apply feature selection during classification. For PKS, later-layer FFNs  
 82 exhibit higher scores for hallucinated responses (Figure 3a), and Pearson correlations show positive  
 83 associations with hallucinations (Figure 3b). Together, these results indicate that RAG hallucinations  
 84 arise from under-utilization of external context and over-reliance on parametric knowledge.

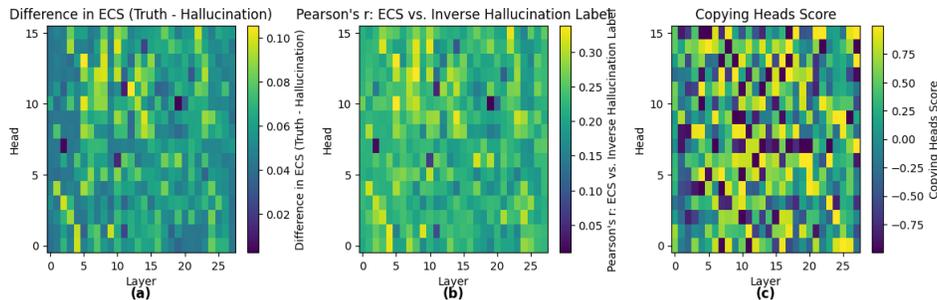


Figure 2: Relationship Between LLM Utilization of External Context and Hallucination

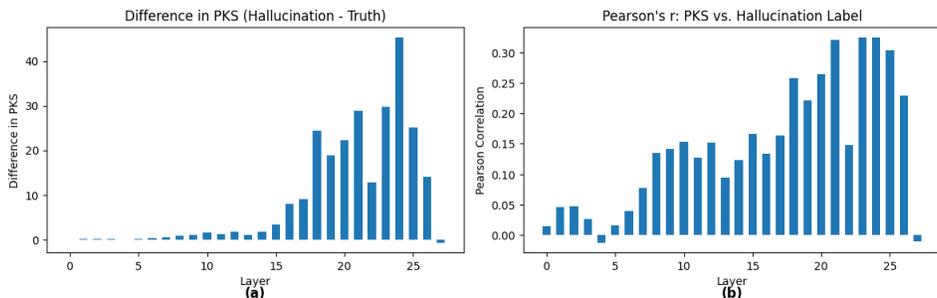


Figure 3: Relationship Between LLM Utilization of Parametric Knowledge and Hallucination

85 **Hallucination Detection** Response-level labels are obtained by aggregating span-level predictions  
 86 from the trained classifier: a response is labeled hallucinated if any span is predicted as such. We

87 evaluate **self-evaluation** (responses from Qwen3-0.6b) and **proxy-based evaluation** (responses  
 88 from other models, e.g., GPT-4.1 mini). We compare against proprietary models (GPT-5, GPT-4.1),  
 89 open-source LLMs (GPT-OSS-20b, LLaMA-3.3-70b, LLaMA-3.1-8b, Qwen3-32b, Qwen3-0.6b),  
 90 and commercial detection systems (RAGAS [16], TruLens [17], RefChecker [18]). Baseline prompts  
 91 and implementation details of commercial detection systems are in Appendices D and E. Results are  
 92 summarized in Table 1. Key observations are:

- 93 • Overall, our method achieves moderate performance. In self-evaluation, it outperforms TruLens  
 94 and LLaMA-3.1-8b and matches RefChecker. In proxy evaluation, it surpasses nearly all models  
 95 except GPT-5 and RAGAS, despite using only a 0.6b parameter model, while higher-performing  
 96 models are proprietary or much larger.
- 97 • We also include Qwen3-0.6b as a detection model to demonstrate that, by itself, the 0.6b model  
 98 is insufficient for hallucination detection. However, when combined with a strong classifier that  
 99 leverages its internal signals, performance is substantially improved.
- 100 • Our model favors recall over precision, likely due to low-confidence false-positive spans retained  
 101 during data curation (Section 2). Nonetheless, we argue that higher recall is generally desirable, as  
 102 false-positive cases can be further verified in downstream.
- 103 • In proxy evaluation, all models (from tiny 0.6b to GPT-5) show higher precision than recall. This  
 104 trend may stem from models such as GPT-4.1-mini producing more reasonable-sounding responses  
 105 where subtle inaccuracies are harder to detect. While in self-evaluation, errors from Qwen3-0.6b  
 106 are more readily detected by stronger models.

Table 1: Response-level Detection Performance (%)

Model	Self-Evaluation			Proxy-based Evaluation		
	Precision	Recall	F1	Precision	Recall	F1
GPT-5	77.27	<b>92.97</b>	<b>84.40</b>	91.67	66.27	<b>76.92</b>
GPT-4.1	76.39	85.94	80.88	<b>94.29</b>	39.76	55.93
GPT-OSS-20b	82.79	78.91	80.80	92.31	43.37	59.02
llama-3.3-70b-versatile	81.03	73.44	77.05	<u>93.75</u>	18.07	30.30
llama-3.1-8b-instant	69.23	49.22	57.53	70.37	22.89	34.55
Qwen3-32b	79.55	82.03	80.77	86.11	37.35	52.10
Qwen3-0.6b	70.27	20.31	31.52	78.79	31.33	44.83
RAGAS	68.45	<u>89.84</u>	77.70	75.29	77.11	<u>76.19</u>
TruLens	<b>89.61</b>	53.91	67.32	49.08	<b>96.39</b>	65.04
RefChecker	<u>84.62</u>	68.75	75.86	71.43	12.05	20.62
Ours	63.89	<u>89.84</u>	74.68	62.90	<u>93.98</u>	75.36

**Note:** RAGAS, TruLens and RefChecker use GPT-4.1 under the hood.

## 107 4 Conclusion and Limitations

108 We developed a RAG hallucinations detection method by decoupling contributions from parametric  
 109 knowledge and external context. Our correlation analysis shows that hallucinations arise from over-  
 110 reliance on parametric knowledge and insufficient use of external context. We trained classifiers to  
 111 predict span-level hallucinations and aggregate them for response-level detection. Additionally, we  
 112 demonstrated our method as a proxy for evaluating large-scale models at low computational cost.

113 Despite these strengths, several limitations remain. Computing ECS and PKS across all layers  
 114 and heads is computationally intensive, particularly when projecting hidden states to vocabulary  
 115 distributions. Future work could identify the most critical layers—our analysis suggests late-layer  
 116 signals may suffice—reducing computation without sacrificing accuracy. We also rely solely on  
 117 ECS and PKS, which limits performance relative to larger models; incorporating uncertainty-based  
 118 or representation-level features could improve detection. Finally, while mechanistic signals could  
 119 enable response intervention to steer models toward truthful outputs, effective real-time manipulation  
 120 would require substantial computational resources and more advanced tooling beyond the current  
 121 proxy-based framework.

## References

- 122
- 123 [1] Siwei Wang et al. Infibench: Evaluating the question-answering capabilities of code large  
124 language models. *OpenReview*, 2024.
- 125 [2] Juyong Jiang et al. A survey on large language models for code generation. *arXiv preprint*  
126 *arXiv:2406.00515*, 2024.
- 127 [3] Zhiwei Ji, Yiming Zhang, Yicheng Liu, et al. A survey on hallucination in large language  
128 models. *arXiv preprint arXiv:2311.05232*, 2023.
- 129 [4] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,  
130 Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models:  
131 A survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023.
- 132 [5] Juntong Song, Xingguang Wang, Juno Zhu, Yuanhao Wu, Xuxin Cheng, Randy Zhong, and  
133 Cheng Niu. Rag-hat: A hallucination-aware tuning pipeline for llm in retrieval-augmented  
134 generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language*  
135 *Processing: Industry Track*, pages 1548–1558, 2024.
- 136 [6] Haichuan Hu, Congqing He, Xiaochen Xie, and Quanjun Zhang. Lrp4rag: Detecting hallucina-  
137 tions in retrieval-augmented generation via layer-wise relevance propagation. *arXiv preprint*  
138 *arXiv:2408.15533*, 2024.
- 139 [7] Likun Tan, Kuan-Wei Huang, and Kevin Wu. Fred: Financial retrieval-enhanced detection and  
140 editing of hallucinations in language models. *arXiv preprint arXiv:2507.20930*, 2025.
- 141 [8] Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song,  
142 and Tong Zhang. Ragtruth: A hallucination corpus for developing trustworthy retrieval-  
143 augmented language models. In *Proceedings of ACL 2024*, 2024.
- 144 [9]  Kov and G Recski. Lettucedetect: A hallucination detection framework for rag  
145 applications. *arXiv preprint arXiv:2501.00232*, 2025.
- 146 [10] Ra V, Timoth Mickus, Elaine Zosa, Teemu Vahtola, J Tiedemann, et al. Semeval-  
147 2025 task 3: Mu-shroom, the multilingual shared task on hallucinations and related observable  
148 overgeneration mistakes. *arXiv preprint arXiv:2504.11975*, 2025.
- 149 [11] Zhongxiang Sun, Xiaoxue Zang, Kai Zheng, Jun Xu, Xiao Zhang, Weijie Yu, and Han Li. Re-  
150 deep: Detecting hallucination in retrieval-augmented generation via mechanistic interpretability.  
151 *arXiv preprint arXiv:2410.11414*, 2024.
- 152 [12] Yuchen Liu, Yiming Zhang, Hao Chen, et al. Lrp4rag: Detecting hallucinations in retrieval-  
153 augmented generation via layer-wise relevance propagation. *arXiv preprint arXiv:2408.15533*,  
154 2024.
- 155 [13] Neel Nanda and Joseph Bloom. Transformerlens. [https://github.com/  
156 TransformerLensOrg/TransformerLens](https://github.com/TransformerLensOrg/TransformerLens), 2022.
- 157 [14] Robert Friel, Masha Belyi, and Atindriyo Sanyal. Ragbench: Explainable benchmark for  
158 retrieval-augmented generation systems. *arXiv preprint arXiv:2407.11005*, 2024.
- 159 [15] Anonymous (TransformerLens Documentation). Transformerlens main demo notebook  
160 – eigenvalue copying score analysis. [https://transformerlensorg.github.io/  
161 TransformerLens/generated/demos/Main\\_Demo.html](https://transformerlensorg.github.io/TransformerLens/generated/demos/Main_Demo.html), 2025. Shows empirical finding  
162 that `OV_copying_score` and `full_OV_copying_score` are “highly (but not perfectly!) correlated”.
- 163 [16] Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. Ragas: Automated  
164 evaluation of retrieval augmented generation. In *Proceedings of the 18th Conference of the*  
165 *European Chapter of the Association for Computational Linguistics: System Demonstrations*,  
166 pages 150–158, 2024.
- 167 [17] Anupam Datta, Matt Fredrikson, Klas Leino, Kaiji Lu, Shayak Sen, Ricardo Shih, and Zifan  
168 Wang. Exploring conceptual soundness with trulens. In *Proceedings of the NeurIPS 2021*  
169 *Competitions and Demonstrations Track*, volume 176, pages 302–307. PMLR, 2022.

170 [18] Xiangkun Hu, Dongyu Ru, Lin Qiu, Qipeng Guo, Tianhang Zhang, Yang Xu, Yun Luo, Pengfei  
 171 Liu, Yue Zhang, and Zheng Zhang. Refchecker: Reference-based fine-grained hallucination  
 172 checker and benchmark for large language models. *arXiv preprint arXiv:2405.14486*, 2024.

## 173 A Preprocessing Example

174 A sample entry from our preprocessed dataset is shown below:

```

175 question: "What is the growth rate in CCG operating income in 2015?"
176
177 documents: [
178   "Management 2019's discussion and analysis of financial condition
179     and results of operations (continued). "
180   "The following results drove changes in CCG operating income by
181     approximately the amounts indicated:"
182 ]
183
184 prompt: "Given the context, please answer the question based on the
185   provided information from the context. Include any reasoning with
186   the answer.\n\nContext: Management 2019's discussion and analysis
187   of financial condition and results of operations (continued). The
188   following results drove changes in CCG operating income by
189   approximately the amounts indicated:\nQuestion: What is the growth
190   rate in CCG operating income in 2015?\nAnswer:"
191
192 response: "The exact growth rate is derived from the differences in
193   the values provided. Thus, the answer is that the growth rate in
194   CCG operating income in 2015 is approximately -20.7%."
195
196 labels: [
197   {start: 673, end: 722, confidence: 0.8037, text: "Growth Rate} = \\  

198     frac{-2161}{10327} \\\times 100\\%"},
199   {start: 730, end: 734, confidence: 0.5833, text: "20.7"}
200 ]
201
202 labels_llama: 1
203 labels_gpt: 1
204
205 prompt_spans: [[0, 135], [136, 144], [145, 348], [349, 1177], [1178, 3465]]
206 response_spans: [[0, 383], [384, 597], [598, 663]]
207
208
  
```

### 209 Notes:

- 210 • labels provide span-level hallucination annotations with start and end indices in the response text.
- 211 • labels\_llama and labels\_gpt indicate whether the respective LLM judges marked the response  
 212 as hallucinated (1) or not (0).
- 213 • prompt\_spans and response\_spans segment the prompt and response into sentence or phrase-  
 214 level chunks for span-level scoring.

## 215 B Computation Details of Mechanistic Metrics

216 **External Context Score:** The External Context Score (ECS) quantifies the extent to which a  
 217 language model leverages external context when generating a response. In mechanistic interpretability,  
 218 attention heads are responsible for retrieving relevant information from the context. To measure this  
 219 utilization, ECS captures the semantic alignment between response segments and the context chunks  
 220 most strongly attended to by the model.

221 Let the external context be partitioned into chunks  $\{\tilde{c}_1, \tilde{c}_2, \dots, \tilde{c}_M\}$ , and the generated output  
 222 into response chunks  $\{\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_N\}$ , corresponding to prompt\_spans and response\_spans in

223 Appendix A. For each attention head  $h$  at layer  $l$ , the most relevant context chunk for each response  
 224 chunk  $\tilde{r}_j$  is identified as

$$\tilde{c}_j^{\ell,h} = \arg \max_{\tilde{c}_i} A(\tilde{r}_j^{\ell,h}, \tilde{c}_i^{\ell,h}),$$

225 where  $A$  denotes the token-level attention weight matrix, and  $l$  and  $h$  indicate the layer and head  
 226 indices. The chunk-level ECS is defined as the cosine similarity between the embeddings of  $\tilde{r}_j$  and  
 227 its corresponding context chunk  $\tilde{c}_j$ :

$$\text{ECS}_{\tilde{r}_j}^{\ell,h} = \cos(e(\tilde{r}_j^{\ell,h}), e(\tilde{c}_j^{\ell,h})),$$

228 where  $e(\cdot)$  represents the embedding function, and  $\cos(\cdot, \cdot)$  denotes cosine similarity.

229 **Parametric Knowledge Score:** The Parametric Knowledge Score (PKS) quantifies the extent to  
 230 which the FFN contributes to *parametric knowledge*, i.e., knowledge stored in the model’s weights,  
 231 as opposed to information coming from external context. This is done by measuring the difference  
 232 between residual stream states **before** the FFN layer and **after** the FFN layer.

233 Since the residual stream itself does not directly indicate "which token is being suggested"—it is  
 234 a latent vector—we map it through the unembedding/projection matrix to the vocabulary distribu-  
 235 tion. This allows us to observe how the residual change after the FFN influences predicted token  
 236 probabilities.

237 We then apply the Jensen-Shannon divergence (JSD) to compute the distance between the two  
 238 vocabulary distributions, which defines the token-level PKS. The chunk-level PKS is computed by  
 239 averaging the token-level PKS over all tokens in the chunk. Mathematically, this is expressed as

$$\text{PKS}_{t_n}^{\ell} = \text{JSD}(p(x_n^{\text{mid},\ell}), p(x_n^{\ell})), \quad \text{PKS}_{\tilde{r}}^{\ell} = \frac{1}{|\tilde{r}|} \sum_{t_n \in \tilde{r}} \text{PKS}_{t_n}^{\ell}.$$

240 where  $p(\cdot)$  denotes the mapping from residual stream states to vocabulary distributions, and  $x_n^{\text{mid},\ell}$   
 241 and  $x_n^{\ell}$  refer to the residual stream states before and after the FFN layer, respectively.  $\text{PKS}_{t_n}^{\ell}$  and  
 242  $\text{PKS}_{\tilde{r}}^{\ell}$  stand for token-level and chunk-level PKS, respectively.

243 The computation was performed using the TransformerLens library on the Qwen3-0.6B model.  
 244 Inference was executed on a Google Colab L4 GPU with 24-GB memory, using `torch.float16`  
 245 precision to reduce activation storage costs. Sentence-level semantic similarity was computed using  
 246 the BAAI/bge-base-en-v1.5 encoder, also hosted on GPU to avoid transfer overhead.

247 The average execution time for an end-to-end ECS/PKS computation was 42 seconds per example.  
 248 GPU allocation remained within 1.9–2.1 GB, with reserved memory peaking at 2.2 GB across  
 249 iterations. After each iteration, tensors were explicitly released, and memory was reclaimed using  
 250 `torch.cuda.empty_cache()` and `torch.cuda.ipc_collect()` to prevent fragmentation.

## 251 C Training Details of Classifiers

252 Training was conducted at the span level using 1,852 instances, corresponding to 7,799 span-level  
 253 samples (4,406 negative, 3,393 positive). We used **External Context Scores (ECS)** and **Parametric**  
 254 **Knowledge Scores (PKS)** as input features. For Qwen3-0.6b (28 layers, 16 attention heads per  
 255 layer), this initially yielded 476 features, which were reduced to 341 after feature selection. The  
 256 dataset was split 90/10 for training and validation with stratification on hallucination labels. Feature  
 257 preprocessing included standardization (`StandardScaler`), removal of near-constant and duplicate  
 258 features (`DropConstantFeatures`, `DropDuplicateFeatures`), and correlation-based filtering  
 259 (`SmartCorrelatedSelection`, Pearson threshold 0.9) using a `RandomForestClassifier` with  
 260 max depth 5. Four classifiers—Logistic Regression, Support Vector Classification (SVC), Random  
 261 Forest, and XGBoost—were evaluated using pipelines combining preprocessing and classification.  
 262 Random Forest and XGBoost had max tree depth 5 with other hyperparameters default. SVC achieved  
 263 the highest validation F1 score and was selected as the final model, while XGBoost, despite strong  
 264 training performance, overfit severely, highlighting the need for regularization or more data. Overall,  
 265 SVC provided the best balance between complexity and generalization (see Table 2).

Table 2: Span-level Detection performance (%)

Classifier	Train Prec.	Val Prec.	Train Rec.	Val Rec.	Train F1	Val F1
LR	79.71	74.84	77.05	71.09	78.36	72.92
SVC	84.44	79.00	79.24	74.34	81.76	76.60
RandomForest	79.87	74.92	76.13	72.27	77.95	73.57
XGBoost	99.74	76.45	99.77	73.75	99.75	75.08

266 **D Prompt for Baselines**

267 Below is the prompt used for hallucination detection when using models GPT-5, GPT-4.1, GPT-OSS-  
 268 20b, LLaMA-3.3-70b-versatile, LLaMA-3.1-8b-instant, Qwen3-32b and Qwen3-0.6b.

```

User Prompt

You are an expert fact-checker. Given a context, a question, and a response,
your task is to determine if the response is faithful to the context.
Context: context
Question: question
Response: response
Is the response supported and grounded in the context above? Answer "Yes"
or "No", and provide a short reason if the answer is "No". Be concise and
objective.
    
```

269

270 **E Implementation of Commercial Detection Systems**

271 We describe our implementation of three commercial tools—RAGAS, TruLens, and RefChecker—for  
 272 hallucination detection below.

273 **RAGAS:** We use RAGAS’s *faithfulness* metric to evaluate how well a model’s responses align  
 274 with the provided context documents. GPT-4.1 is used as the evaluator model. For each data point, a  
 275 faithfulness score between 0 and 1 is computed. To determine the optimal threshold that maximizes  
 276 F1 score, we evaluate F1 across candidate thresholds in [0,1]. Predictions are binarized at each  
 277 threshold, and the threshold that maximizes F1—balancing precision and recall—is selected.

278 **TruLens:** TruLens provides a framework for evaluating hallucination via its groundedness feed-  
 279 back mechanism. We use the `groundedness_measure_with_cot_reasons` function to compute  
 280 groundedness scores in [0,1]. The F1-optimal threshold is determined using the same procedure as in  
 281 RAGAS.

282 **RefChecker:** RefChecker extracts factual claims from model responses and verifies them against  
 283 reference documents. It consists of two components: an `LLMExtractor` that extracts claims, and  
 284 an `LLMChecker` that classifies claims as `Entailment`, `Neutral`, or `Contradictory`. A response is  
 285 labeled hallucinated if any claim is `Contradictory`, and non-hallucinated if all claims are either  
 286 `Entailment` or `Neutral`.