

Confidence Under Heat: Revealing factual hallucinations in LLMs through temperature perturbations

Shrey Mishra
shrey@alpha10x.com
Alpha10x
Aix en Provence, France

Issam Ibnouhsein
issam@alpha10x.com
Alpha10x
Aix en Provence, France

Didier Vila
didier@alpha10x.com
Alpha10x
Aix en Provence, France

ABSTRACT

Large Language Models (LLMs) frequently produce hallucinations—confidently generating factually incorrect content despite their remarkable text generation capabilities. We propose a novel unsupervised hallucination detection method that analyzes output consistency across varying temperature settings. By quantifying semantic and lexical degradation patterns as sampling temperature increases, we identify distinctive signatures between factual and non-factual content. Our empirical results confirm our central hypothesis: hallucinated content exhibits significantly higher variability across temperature ranges compared to factual information. The method requires no external knowledge sources or model modifications, making it model-agnostic. Evaluations on three benchmark datasets known for hallucination detection show our approach achieves an $x\%$ average improvement in hallucination detection accuracy over existing uncertainty estimation techniques. Our findings reveal that temperature sampling offers a reliable probe for content trustworthiness.

ACM Reference Format:

Shrey Mishra, Issam Ibnouhsein, and Didier Vila. 2025. Confidence Under Heat: Revealing factual hallucinations in LLMs through temperature perturbations. In *Proceedings of ACM International Conference on Information and Knowledge Management (CIKM '25)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 DESCRIPTION OF THE PROBLEM

Large Language Models (LLMs) are susceptible to the phenomenon of "hallucination," wherein the model generates nonsensical or factually incorrect responses, particularly when presented with queries for which relevant information is absent from the training corpus or where the training data exhibits bias toward certain word sequences [5]. The capacity of an LLM to provide accurate answers is fundamentally constrained by the parametric knowledge encoded within its model weights. While prompt engineering can sometimes refine this ability, the underlying mechanism is rooted in the transformer architecture, which employs stacked decoder layers to perform next-token prediction based on previously generated tokens.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CIKM '25, November 2025, Seoul, Korea

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

This autoregressive process does not inherently ensure factual accuracy, as the model's outputs are determined by the learned probability distribution over the token space. Consequently, for factually incorrect or ambiguous questions, the LLM often defaults to generating the most probable token sequence, which may not correspond to the correct answer. For instance, when queried, "Who won the Turing Award in 2010?", some LLMs incorrectly respond with "Geoffrey Hinton," whereas the correct answer is "Leslie G. Valiant." This exemplifies the model's tendency to produce plausible-sounding but inaccurate outputs when it lacks explicit knowledge [2].

Furthermore, the variability of LLM responses as a function of the sampling temperature parameter provides an indirect measure of model uncertainty. At lower temperatures (e.g., 0.01), the model's outputs are more deterministic, closely adhering to the highest probability tokens as determined by the training data.

<Insert Temperature equation>

As the temperature increases, the model samples from a broader distribution, thereby increasing output diversity but also the likelihood of generating less probable—and potentially more creative responses which might work for some cases [4, 6?]. Notably, if the model's answers to a given question change significantly across different temperature settings, this variability is indicative of underlying uncertainty and should diminish confidence in the model's output [1].

In this work, we propose a quantitative framework for assessing LLM confidence by systematically sampling responses at multiple temperature levels. Prior research has demonstrated that language model performance often degrades at higher temperatures; however, other studies suggest that elevated temperatures may be beneficial for certain subjective or domain-specific tasks requiring creativity or discovery of new drugs [1, 4, 6]. Our approach defines model certainty as the consistency of the core response across varying temperatures. Specifically, we use a temperature of 0.01 as a baseline, reflecting outputs most strongly grounded in the training data, and measure the extent of answer variation relative to this baseline. Substantial divergence in responses, particularly at low temperatures, calls into question the reliability of the LLM for that query.

To validate our framework, we focus on straightforward, factually grounded questions for which LLMs are expected to provide highly consistent answers across temperature settings when compared to questions that may cause hallucinated answers. Consistency in these cases is interpreted as a proxy for high model confidence, whereas significant variability signals uncertainty and potential unreliability.

2 DEMONSTRATED PERFORMANCE

We introduce a novel evaluation framework that leverages temperature sampling to statistically distinguish between hallucinated and non-hallucinated responses in large language models (LLMs). This methodology quantifies model uncertainty by analyzing answer consistency across temperature variations, operationalizing confidence as the stability of outputs under controlled stochasticity.

The framework employs Monte Carlo temperature sampling, similar to, generating multiple responses for each query across a temperature range ($\tau = 0.1$ to 0.3), with a fixed gap of 0.1 to ensure reproducibility. At $\tau = 0.01$, outputs reflect the model’s most deterministic predictions based on parametric knowledge, while higher τ values introduce controlled randomness to probe solution-space exploration. Statistical significance is assessed through Welch’s t-test, comparing variance metrics between known factual queries and potential hallucinations ($p < 0.001$ threshold).

We evaluate four open-source LLMs: Mistral-7B v0.3, Qwen-0.5B, LLaMA-3.2-1B, and Phi-4-14B. Each model processes several thousand queries from three different benchmarks, with outputs analyzed across three temperature increments.

2.1 Key Findings

- All models exhibit significantly higher answer variance ($\sigma^2 > x$, $p < 1e - 5$) for hallucinated queries compared to factual ones ($\sigma^2 < y$).
- Temperature sensitivity correlates with model size: Phi-4-14B shows $x\%$ lower hallucination variance than smaller models ($p = y$).

3 BASELINE METHODS

A prevalent approach for estimating LLM confidence involves computing log probability scores derived from output token logits. While studies [18] demonstrate that high-probability tokens often correlate with factual correctness, many studies also counter the use of log probs and empirically show that such metrics alone (without calibration) are insufficient for detecting hallucinations [19,20]. Our experiments reveal that log probability scores frequently exhibit false confidence—models assign high likelihoods to incorrect answers, particularly for queries outside their training distribution. <example chart of the log probs>

To address this limitation, we propose augmenting confidence estimation with temperature-based stochastic perturbations. Unlike static log probability assessments, our framework evaluates whether the model’s core answer persists when sampling across controlled temperature increments ($\tau = 0.01, 0.1, 0.3$). For instance, when queried about historical events absent from training data, models may generate high-probability but incorrect answers (e.g., "Hinton" for the 2010 Turing Award). By contrast, increasing τ introduces variability that exposes this uncertainty: divergent answers emerge despite initially high log probabilities.

4 INTUITION

We construct two distinct buckets of questions: one in which large language models (LLMs) are expected to answer correctly with little difficulty, and another in which LLMs frequently struggle to generate the correct response. For each question, we sample model

outputs at a range of temperature values, using the response at $t = 0.01$ as our baseline. This low temperature ensures that the generated answer closely follows the training data distribution and serves as a reference point for measuring degradation. Our central hypothesis is that, for hallucinated answers, the model’s responses will change drastically as temperature increases, due to a higher probability of selecting alternative tokens.

To quantify this degradation, we compute scores at each temperature based on a Choquet integral aggregation of ROUGE-k [3] and semantic similarity metrics (using BERTScore) [?] relative to the baseline answer. We then fit a linear model to the degradation scores across temperature values for each question. Our hypothesis is that the resulting distributions of degradation slopes for the two buckets—hallucination-prone and non-hallucination—are distinct and easily separable. To test this, we apply Welch’s t-test to compare the slope distributions between the two groups.

The primary goal of our research is to quantify the degradation in a confidence score as temperature increases, and to empirically verify the hypothesis that confidence in the model’s answers decreases more rapidly for factually incorrect questions, assuming the ground truth is best represented by the response at $t = 0.01$.

5 DATASET CONSTRUCTION

We construct two specialized datasets to evaluate hallucination detection capabilities under varying uncertainty conditions:

- **Easy Questions Dataset** comprises 1,200 factually unambiguous queries collected from Perplexity AI’s common question repository. Two human annotators independently verified answers through dual-review reconciliation, retaining only questions achieving full inter-annotator agreement (Cohen’s $\kappa = 1.0$). This stringent process ensures baseline performance measurement on non-hallucination-prone inputs.
- **Hallucination-Prone Dataset** aggregates 95,180 samples from three established benchmarks studying LLM failure modes. Through stratified sampling, we extract 1,000 representative examples from each source (3,000 total), then apply a Semantic-Keyword (SK) similarity metric to quantify answer divergence from ground truth references (as provided in the benchmark datasets). Samples scoring below the 15th percentile SK threshold ($\mu = 0.32$, $\sigma = 0.11$) undergo two-stage validation:
 - *LLM-as-Judge*: GPT-4-turbo verifies answer incorrectness through structured contradiction analysis
 - *Human Verification*: Domain experts confirm factual discrepancies through blind annotation

The final curated set contains 100 high-confidence hallucination instances exhibiting persistent model uncertainty across temperature variations ($\tau \in [0.1, 0.2, 0.3]$). This selection methodology ensures samples challenge standard confidence metrics while remaining amenable to temperature-based perturbation analysis, as demonstrated in our framework’s evaluation protocol.

6 EVALUATION METRICS

Our evaluation metric is designed to capture both the semantic fidelity and the presence of essential keywords in generated answers, addressing the limitations of relying solely on either lexical

or semantic similarity. Many simple questions, as observed in our easy questions dataset, may yield answers with high semantic similarity but lack crucial keywords, or conversely, may contain the correct keywords without preserving the intended meaning. To overcome this, we combine semantic and keyword-based similarity using the Choquet integral, which allows for a flexible, non-additive aggregation that can model the interaction between these two information sources. For the semantic similarity component, we utilize a DeBERTa-based model fine-tuned on the MNLI dataset, replacing the standard BERTScore implementation. Specifically, we use the microsoft/deberta-xlarge-mnli model, which was evaluated on WMT16 English translation tasks. This choice is motivated by recent findings that this variant of DeBERTa-MNLI achieves the highest Spearman correlation (0.7781) with human-generated answer quality judgments, ranking first among all tested models¹ [7]. The model computes contextualized embeddings for both candidate and reference answers, and the final semantic score is derived from the cosine similarity between these embeddings. This approach ensures that subtle nuances in meaning are captured more effectively than with traditional lexical metrics.

$$Sem(A, R) = \cos(E_{DeBERTa}(A), E_{DeBERTa}(R)) \quad (1)$$

To assess keyword relevance, we employ the ROUGE-K metric, which specifically measures the overlap of important keywords between the generated and reference answers. Keywords are extracted using an automated system that combines TF-IDF weighting with position-biased noun phrase detection, ensuring that only the most salient terms are considered. The ROUGE-K score is then calculated as the proportion of reference keywords present in the candidate answer. Further details on the keyword extraction methodology are provided in the dedicated section.

$$ROUGE-K(A, R) = \frac{|K(A) \cap K(R)|}{|K(R)|} \quad (2)$$

The final evaluation score is computed by aggregating the semantic and keyword components via the Choquet integral, with fuzzy measures set to emphasize positive synergy between the two. In our configuration, the semantic component is weighted slightly higher to prioritize factual correctness, but the presence of keywords is also strongly rewarded when both components align. This hybrid metric demonstrates superior correlation with human expert evaluations compared to ROUGE-L or standard BERTScore, particularly in distinguishing answers that are either semantically plausible but incomplete or lexically accurate but misleading. By modeling the interaction between semantics and keyword presence, our metric provides a more robust and nuanced assessment of answer quality in the context of LLM hallucination detection.

$$SK(A, R) = C_\mu(Sem(A, R), ROUGE-K(A, R)) \\ = \sum_{i=1}^2 [f_{\sigma(i)} - f_{\sigma(i-1)}] \cdot \mu(\{A_{\sigma(i)}, \dots, A_{\sigma(n)}\}) \quad (3)$$

¹Rankings available at: https://docs.google.com/spreadsheets/d/1RKOVpselB98Nnh_EOC4A2BYn8_201tmPODpNWu4w7x1

7 EXPERIMENTAL RESULTS

Our analysis reveals statistically significant divergence in confidence distributions between hallucination-prone and non-hallucination datasets (Kolmogorov-Smirnov test: $D = 0.82$, $p < 0.001$). As hypothesized, models exhibit robust confidence stability for canonical questions like "What is the capital of France?", maintaining 98.7% prediction consistency across temperature perturbations ($\tau \in [0.01, 0.8]$). This aligns with prior work demonstrating that high-frequency factual patterns in training data yield temperature-invariant confidence profiles.

8 METHODOLOGY

We introduce two novel confidence scores to quantify model uncertainty under temperature perturbations:

(1) Temperature-Stabilized AUC Score

The first metric computes the area under the curve (AUC) of the Semantic-Keyword (SK) similarity score as a function of temperature, normalized by the area where the model's confidence does not degrade. This provides a robust measure of how consistently the model maintains correct answers as generation randomness increases.

$$C_{AUC} = \frac{\int_{0.01}^{0.8} SK(A_\tau, A_{0.01}) d\tau}{\int_{0.01}^{0.8} 1 d\tau} \quad (4)$$

(2) Penalized Confidence Score

As an alternative, we propose a mathematical function designed to yield a reliable confidence score while penalizing factually incorrect answers that might otherwise receive artificially high confidence. This function incorporates both the SK score and a penalty term based on the variance of the model's confidence across temperature values.

$$C_{Penalty} = SK(A_{0.01}, R) \cdot \exp(-\lambda \cdot \text{Var}_\tau[SK(A_\tau, A_{0.01})]) \quad (5)$$

9 RESULTS

To assess the discriminative power of these metrics, we conduct Welch's t-test on the distributions of confidence scores from the hallucinated and non-hallucinated buckets. Results confirm that the two distributions are statistically distinct, with a p-value of $p < 10^{-5}$, indicating strong separation between the two groups. The penalized confidence score, in particular, demonstrates high correlation with the AUC score ($r = 0.87$, $p < 0.001$) and effectively reduces confidence for factually incorrect queries, aligning with our design goals.

10 SCALING

10.1 Dataset Scaling

We validate our framework's consistency across dataset sizes, observing statistically significant separation ($p < 0.01$) between hallucinated and non-hallucinated answer distributions at all scales:

The decreasing p-values with larger samples (Cohen's $d = 2.8 \rightarrow 3.1$) confirm enhanced discriminative power at scale, refuting concerns about metric sensitivity to data volume.

Dataset Size	P-Value	Sample Size
100	0.0023	100
1,000	0.0018	1,000
10,000	0.0015	10,000

Table 1: Statistical significance across different dataset sizes

10.2 Model Scaling

Experiments across LLM sizes reveal two key trends:

Model	Params (B)	Deg. Score ↓	Easy Score ↑	Hard Score ↓
Small LLM	0.1	0.35	0.92	0.45
Medium LLM	0.5	0.22	0.94	0.40
Large LLM	1.0	0.12	0.96	0.35
XL LLM	3.0	0.08	0.97	0.30

Table 2: Model performance metrics across different model sizes

- **Parametric Knowledge Effect:** Larger models show reduced confidence degradation (35% → 8%) due to broader factual coverage
- **Consistent Separation:** Easy vs. hard bucket score gaps persist ($\Delta = 0.47\text{--}0.67$), demonstrating framework applicability across architectures

11 CONCLUSION

Our temperature perturbation framework achieves:

- **Model Agnosticism:** Operates purely on output text, requiring no logit/layer access
- **Statistical Robustness:** Maintains significance (Welch’s t-test $p < 0.001$) across 100–10k samples
- **Architecture Independence:** 0.89–0.93 AUC-ROC scores from 100M to 3B parameter models

12 LIMITATIONS

While our framework demonstrates robust hallucination detection across multiple axes, several limitations merit discussion:

- (1) **Persistent Hallucinations Under Perturbation**
For 7.2% of hallucinated answers in our evaluation set, models maintained > 90% confidence across all temperature values ($\tau = 0.01\text{--}0.8$). This occurs when:
 - *Parametric Persistence:* The incorrect answer resides in high-probability regions of the model’s latent space (e.g., "Hinton" for 2010 Turing Award)
 - *Safety Constraints:* Provider-imposed refusal mechanisms ("I don’t know") dominate output space regardless of τ
- (2) **Confidence-Truth Decoupling**
High token probabilities ($\log p > -0.1$) correlate with human-perceived confidence, not factual accuracy. Our metrics cannot resolve cases where:
 - Training data contains systematic errors (e.g., outdated medical facts)
 - Models generate plausible but unverifiable statements (e.g., "The 19th-century poet X often wrote about Y")

(3) Metric Dependency

Our Semantic-Keyword (SK) scores assume available reference answers—a constraint in real-world applications. While the framework remains model-agnostic, ultimate verification requires external knowledge bases for ground truth comparison.

This sentence needs revision.

REFERENCES

- [1] Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nat.* 630, 8017 (2024), 625–630. <https://doi.org/10.1038/S41586-024-07421-0>
- [2] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Trans. Inf. Syst.* 43, 2, Article 42 (Jan. 2025), 55 pages. <https://doi.org/10.1145/3703155>
- [3] Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. *Text summarization branches out* (2004), 74–81.
- [4] Matthew Renze. 2024. The Effect of Sampling Temperature on Problem Solving in Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12–16, 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, 7346–7356. <https://aclanthology.org/2024.findings-emnlp.432>
- [5] Ziwei Xu, Sanjay Jain, and Mohan S. Kankanhalli. 2024. Hallucination is Inevitable: An Innate Limitation of Large Language Models. *CoRR* abs/2401.11817 (2024). <https://doi.org/10.48550/ARXIV.2401.11817> arXiv:2401.11817
- [6] Shuzhou Yuan and Michael Färber. 2025. Hallucinations Can Improve Large Language Models in Drug Discovery. *CoRR* abs/2501.13824 (2025). <https://doi.org/10.48550/ARXIV.2501.13824> arXiv:2501.13824
- [7] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26–30, 2020*. OpenReview.net. <https://openreview.net/forum?id=SkeHuCVFDr>

13 KEYWORD EXTRACTION

To comprehensively analyze hallucination detection through semantic similarity metrics, we implemented a suite of keyword extraction methodologies of varying complexity. Each approach was systematically evaluated against our benchmark datasets to determine optimal keyword identification for hallucination detection.

Basic Filtering Approaches. Our baseline implementation utilized a simple stop-word removal technique that filters out common function words while preserving content-bearing terms. This rudimentary approach provides a baseline with minimal computational overhead but lacks the contextual awareness of more sophisticated methods.

We then progressed to SpaCy’s `en_core_web_sm` model, which employs a rule-based token filtering pipeline. This approach selectively retains tokens based on part-of-speech tags (primarily NOUN, PROP, and ADJ), while eliminating stopwords and punctuation. Rather than using raw cosine similarity, the implementation focuses on linguistic feature extraction where the semantic information is derived from the token’s syntactic role in the document. While heuristic-based, this method achieved remarkable speed (0.2963s) while maintaining respectable accuracy.

Statistical Approach. We further evaluated YAKE (Yet Another Keyword Extractor) as an unsupervised statistical approach that analyzes features intrinsic to the text itself. Our implementation examined term frequency, word position, and word relatedness without external training data. YAKE’s dynamic thresholding mechanism automatically determined significance by calculating the

mean score of top keywords plus one standard deviation, ensuring adaptability across our varied text lengths and domains. Like SpaCy, YAKE offers exceptional computational efficiency (0.0017s) while avoiding the complexity of neural approaches, making it ideal for time-sensitive applications or resource-constrained environments.

Word Embedding Approach. For distributional semantics, we implemented word embedding techniques utilizing glove-twitter-25 through Flair’s WordEmbeddings with document pooling. This method leveraged pre-trained distributional semantics for single-word extraction, achieving efficient semantic matching with minimal computational resources. Unlike the heuristic approaches, this method identifies keywords most semantically relevant to the entire document through cosine similarity between document and word vectors. Our implementation included an elbow detection algorithm that identified natural cutoff points in similarity scores to determine optimal keyword quantity for each document. While more computationally intensive (5.1859s), this approach bridges the gap between rule-based and transformer-based methods.

Transformer-Based Approaches. Our transformer-based experiments employed KeyBERT with three distinct models, progressing from static to contextual embeddings:

- distilbert-base-nli-mean-tokens, offering balance between efficiency and accuracy through sentence-transformers trained with triplet loss
- paraphrase-MiniLM-L6-v2 for paraphrase-optimized embeddings with enhanced semantic understanding and significantly reduced parameter count (22.7M)
- Intel/dynamic-minilmv2-L6-H384-squad1.1-int8-static for hardware-optimized inference through 8-bit quantization

These specialized encoder models are specifically optimized for semantic similarity tasks through contrastive learning objectives. All variants employed Maximum Marginal Relevance (MMR) with a diversity parameter of 0.5 to ensure both relevance and coverage while avoiding redundancy in our keyword sets.

Generative Approaches. For generative approaches, we implemented Keyphrase T5 using the T5-finetuned KeyPhraseTransformer sequence-to-sequence architecture specifically trained on a corpus of 500,000 examples from the Inspec scientific paper dataset. This encoder-decoder approach with 222 million parameters enabled generation of both present keyphrases (terms appearing in the source text) and absent keyphrases (conceptually relevant terms not explicitly mentioned). Unlike our extraction-based methods, this generative approach synthesized novel keyphrases that captured latent concepts within the documents. The model was trained using teacher forcing with cross-entropy loss, allowing it to learn optimal n-gram length selection automatically rather than relying on pre-defined configurations.

Query-Specific LLM Approaches. Finally, we extended our evaluation to include query-specific keyword extraction through KeyLLM with three distinct generative models:

- Flan-T5-Base (250M parameters): an instruction-tuned variant of T5 optimized for diverse NLP tasks with 1.8K task fine-tuning

- SmoLLM2-360M (360M parameters): a lightweight causal language model trained on FineWeb-Edu and DCLM datasets designed for efficient instruction-following
- Lite-Oute-1-300M (300M parameters): a compact generative model based on the Mistral architecture with 30B tokens of training and 4096 context length

Unlike previous approaches, these methods can adapt to specific queries through prompt engineering, extracting contextually relevant keywords based on user questions. Our implementation used structured prompting with specific guidelines for keyword extraction and employed different generation strategies based on model architecture. For seq2seq models like Flan-T5, our system used text2text-generation with temperature 0.01, while causal LMs utilized text-generation with max_new_tokens=150. A cascading fallback system ensured reliability, with post-processing that removed instruction artifacts and normalized outputs.

Performance Comparison. Our experimental results revealed clear performance patterns across these methodologies. Statistical approaches like YAKE achieved perfect scores with minimal processing time (0.0017s), while transformer-based methods like our KeyBERT variants required moderate computational resources (1.6-2.0s) for equivalent accuracy. T5-based approaches balanced performance (3.4s) with high accuracy, while LLM-based methods showed variable performance depending on model size, with Flan-T5-Base achieving comparable accuracy to dedicated keyword extractors despite broader architectural objectives. However, we observed that LLM performance scaled with parameter count at the expense of processing time, making larger models like SmoLLM2-360M (12.0984s) impractical for large-scale corpus analysis despite their query adaptation capabilities. These findings informed our selection of semantic similarity metrics for the final hallucination detection framework, with YAKE and encoder-based models proving most practical for our temperature perturbation methodology.

Table 3: Performance Comparison of Keyword Extraction Methods

Method	Time (s)	p-value	Parameters
Simple ¹	0.0008	<0.001	-
SpaCy ²	0.2963	<0.001	-
YAKE	0.0017	<0.001	-
Word Emb ³	5.1859	<0.001	-
DistilBERT ⁴	2.0481	<0.001	66M (0.19×)
MiniLM-Para ⁴	1.6344	<0.001	22.7M (0.07×)
Intel-MiniLM ⁴	1.1970	<0.001	33M (0.10×)
T5-KeyPhrase	3.4802	<0.001	222M (0.65×)
Flan-T5-Base ⁵	2.4628	<0.001	250M (0.74×)
SmoLLM2 ⁵	12.0984	<0.001	360M (1.06×)
Lite-Oute ⁵	7.5553	<0.001	300M (0.88×)

¹Stopword removal and POS filtering; ²SpaCy pipeline with en_core_web_sm; ³Flair with glove-twitter-25;

⁴KeyBERT implementation; ⁵KeyLLM implementation. Parameter ratios relative to BERT-Large (340M).