# Beyond Semantic Entropy: Boosting LLM Uncertainty Quantification with Pairwise Semantic Similarity

**Dang Nguyen**
UCLA CS
dangnth@cs.ucla.edu

**Ali Payani**
Cisco Systems Inc.
apayani@cisco.com

**Baharan Mirzasoleiman**
UCLA CS
baharan@cs.ucla.edu

## Abstract

Hallucination in large language models (LLMs) can be detected by assessing the uncertainty of model outputs, typically measured using entropy. Semantic entropy (SE) enhances traditional entropy estimation by quantifying uncertainty at the semantic cluster level. However, as modern LLMs generate longer one-sentence responses, SE becomes less effective because it overlooks two crucial factors: intra-cluster similarity (the spread within a cluster) and inter-cluster similarity (the distance between clusters). To address these limitations, we propose a simple black-box uncertainty quantification method inspired by nearest neighbor estimates of entropy. Our approach can also be easily extended to white-box settings by incorporating token probabilities. Additionally, we provide theoretical results showing that our method generalizes semantic entropy. Extensive empirical results demonstrate its effectiveness compared to semantic entropy across two recent LLMs (Phi3 and Llama3) and three common text generation tasks: question answering, text summarization, and machine translation. Our code is available at https://github.com/BigML-CS-UCLA/SNNE.

## 1 Introduction

Large Language Models (LLMs) have demonstrated impressive capabilities in understanding and generating human-like text, revolutionizing various fields (Bubeck et al., 2023; Team et al., 2024a). However, they are not without flaws, and one of the most significant challenges is hallucination, i.e., incorrect or fabricated information that appears plausible (Maynez et al., 2020; Ji et al., 2023). Detecting and mitigating hallucination is critical to ensuring the reliability and safety of LLMs. A common approach to addressing this issue is uncertainty quantification (UQ) (Huang et al., 2024). By measuring the uncertainty of an LLM's outputs, we can identify potentially hallucinated content and flag it for further review, improving both accuracy and user trust.

A straightforward way to estimate uncertainty in LLMs is to generate multiple responses and aggregate their token-level likelihoods (Malinin and Gales, 2020). However, this approach ignores semantic information, treating reworded yet equivalent answers as distinct. Semantic Entropy (SE) (Farquhar et al., 2024) addresses this by clustering semantically similar outputs using bidirectional entailment predictions from a NLI model before computing entropy. While effective for short responses, SE struggles when model generates long one-sentence outputs, a pattern inherent to tasks such as summarization (Zhang et al., 2020) and translation (Hendy et al., 2023). In these cases, SE often reverts to the previous naive approach while its discrete version, Discrete Semantic Entropy (DSE), yields constant values. These shortcomings highlight the need for more robust uncertainty estimation techniques tailored to complexities of lengthy model responses.

We empirically observe that state-of-the-art LLMs tend to generate longer one-sentence outputs on common QA benchmarks: Llama-3.1-8B and Phi-3-mini produce responses with an average length of 4.1 and 4.9 words, respectively, compared to 2.3 words for Llama2-7B. In these scenarios, model outputs are more likely to belong to different semantic groups, thus capturing inter-cluster similarity is important. In addition, each semantic cluster has different spread, necessitating the consideration of inter-cluster similarity. To address these two issues, we propose a simple yet effective black-box UQ method inspired by nearest neighbor estimates of entropy (Beirlant et al., 1997). Our approach can be seamlessly extended to white-box settings by incorporating token probabilities, providing flexibility across different use cases. Additionally, we theoretically prove that our method generalizes semantic entropy, offer-
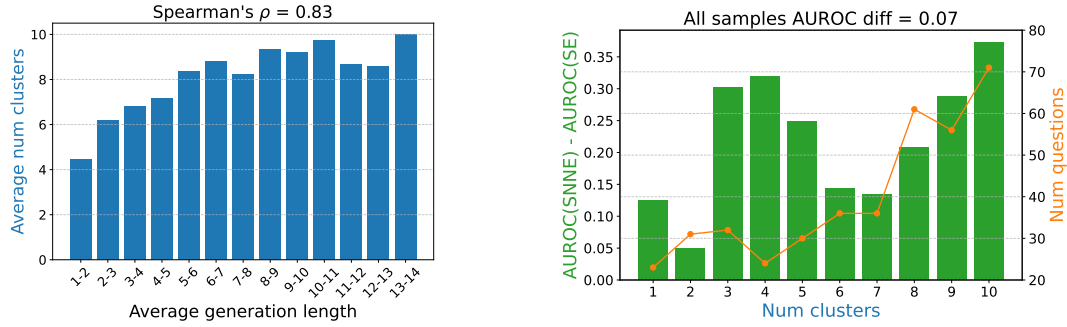
Figure 1: **Generated answers of Llama-3.1-8B on SQuAD.** (left) The Spearman's rank correlation coefficient between the average output length and the number of semantic clusters is 0.83, indicating a high correlation. (right) The AUROC difference between SNNE and SE when calculating on subsets of questions with different number of clusters. SNNE is consistently more distinctive than SE, especially when the number of clusters is large.

ing greater expressiveness. Extensive experiments on nine datasets and state-of-the-art LLMs demonstrate that our method consistently outperforms SE and other baselines across question answering, summarization, and translation.

## 2 Existing uncertainty quantification methods

**Black-box methods.** These methods rely solely on generated answers. pTrue (Kadavath et al., 2022) estimates uncertainty by prompting the model a True/False question to verify its own responses. The model's uncertainty is then determined by the probability of selecting False. Lexical Similarity (LexSim) (Fomicheva et al., 2020) measures uncertainty via average lexical overlap among responses but struggles with semantically similar outputs that use different wording. Graph-based approaches (Lin et al., 2024) structure responses as a similarity graph and derive uncertainty from its properties. NumSet counts semantic clusters found by SE, while SumEigv generalizes this via the sum of eigenvalues of the graph Laplacian. Deg computes the average pairwise semantic similarities, and Eccen quantifies uncertainty based on distances from responses to their center in the space defined by the $k$ smallest eigenvectors of the graph Laplacian. Because graph-based methods use summation to aggregate answer similarities, they are sensitive to outliers or peculiar answers. LUQ (Zhang et al., 2024) is designed for multi-sentence scenarios by aggregating atomic uncertainty scores computed for individual sentences. However, its atomic scoring relies on NLI-based similarity, which inherits the limitations of LexSim and graph-based methods. Our method can be integrated into LUQ to provide more reliable atomic scores, offering a promising direction for

extending our approach to multi-sentence settings. We leave this for future work.

**White-box methods.** These methods utilize internal model information, such as token likelihoods and representations, to assess uncertainty. Naive Entropy (NE) (Malinin and Gales, 2020) computes entropy over length-normalized sequence probabilities but is sensitive to variations in probability assignments for semantic equivalence answers, leading to unreliable estimates. Semantic Entropy (SE) (Farquhar et al., 2024) improves robustness by leveraging a greedy bi-directional entailment clustering algorithm to group responses into semantic clusters before computing entropy but does not account for semantic distances between generations. Kernel Language Entropy (KLE) (Nikitin et al., 2024) addresses this by using von Neumann entropy on a semantic kernel. However, the use of heat and Matern kernels for constructing semantic kernels in KLE makes it challenging to interpret the model's uncertainty and may lead to the loss of information. Empirically, we show that directly utilizing semantic similarities between answers leads to improved performance and more accurate uncertainty estimation. In addition, its $\mathcal{O}(N^3)$ complexity makes it impractical for large-scale uncertainty estimation. In contrast, SE provides a more computationally efficient approach $\mathcal{O}(N^2)$ but struggles with long outputs. Given the efficiency and intuitive formulation, SE serves as a foundation for further improvement, motivating our approach to address its limitations.

Rather than relying on clustering, SAR (Duan et al., 2023) computes a soft aggregation of word-or sentence-level probabilities weighted by their semantic similarity. It also reduces the impact of irrelevant tokens and low-quality sequence samples. Although SAR achieved strong performance

on the LM-Polygraph benchmark (Vashurin et al., 2025), its effectiveness degrades with increasing output length, as shown in their Table 4. Eigen-Score (Chen et al., 2024), in contrast, assumes deeper access to the model's internal states. It leverages the eigenvalues of the covariance matrix of output embeddings—a more computationally intensive approach compared to methods based on logits or probabilities. Additionally, EigenScore applies feature clipping to suppress extreme values in the embedding space, aiming to avoid overconfident, self-consistent hallucinations. However, this technique requires tuning a clipping threshold, which can introduce additional complexity.

## 3 Preliminaries

**Notations.** Let $P(\boldsymbol{a}|\boldsymbol{q})$ denote the sequence probability that the model generates the answer $\boldsymbol{a}$ given question $\boldsymbol{q}$, i.e, $P(\boldsymbol{a}|\boldsymbol{q}) = \sum_j \log p(\boldsymbol{a}_j|\boldsymbol{q} \oplus \boldsymbol{a}_{<j})$. The length-normalized sequence probability is defined as $\tilde{P}(\boldsymbol{a}|\boldsymbol{q}) = P(\boldsymbol{a}|\boldsymbol{q})/len(\boldsymbol{a})$.

**Uncertainty quantification pipeline.** In general, each UQ method consists of two main steps.

1. Generation: Given question $\boldsymbol{q}$, sample $n$ answers $\boldsymbol{a}^1, \ldots, \boldsymbol{a}^n$ from LLM.

2. Entropy estimation: Compute uncertainty based on question $\boldsymbol{q}$ and answers $\{\boldsymbol{a}^i\}_{i=1}^n$.

**Semantic Entropy.** SE (Farquhar et al., 2024) uses a bidirectional entailment model to cluster outputs into $M$ semantic classes $\{C_k\}_{k=1}^M$. Let the semantic class probability be the sum of sequence probabilities of all outputs in that class, i.e., $P(C_k) = \sum_{i,\boldsymbol{a}^i \in C_k} \tilde{P}(\boldsymbol{a}^i|\boldsymbol{q})$. Then, define the normalized semantic class probability as $\bar{P}(C_k) = \frac{P(C_k)}{\sum_{j=1}^M P(C_j)}$. Using the Rao-Blackwellized Monte Carlo estimator, SE computes the entropy as

$$\text{SE}(\boldsymbol{q}) = -\sum_{k=1}^M \bar{P}(C_k) \log \bar{P}(C_k) \qquad (1)$$

**Discrete Semantic Entropy.** DSE (Farquhar et al., 2024) is an extension of SE to the black-box setting by approximating the normalized semantic class probability with the empirical cluster probability, i.e., $\frac{|C_k|}{n}$. Then, the formula of DSE reads

$$\text{DSE}(\boldsymbol{q}) = -\sum_{i=1}^M \frac{|C_k|}{n} \log \frac{|C_k|}{n} \qquad (2)$$

**Issues of Semantic Entropy.** Empirically, we observe that state-of-the-art LLMs tend to generate longer responses than their predecessors. For instance, on five QA datasets in our experiments, Llama-3.1-8B and Phi-3-mini produce responses with an average length of 4.1 and 4.9 words, respectively, compared to 2.3 words for Llama2-7B. Additionally, as shown in Figure 1 left, the number of semantic clusters is strongly correlated with response length. When the number of semantic clusters $M$ approaches $n$, DSE produces a constant entropy regardless of $\{\boldsymbol{a}^i\}_{i=1}^n$ because it fails to account for the similarity among clusters (inter-cluster similarity). Even if the number of clusters $M$ is small in Figure 1 right, SE falls behind our method as it lacks the spread of different clusters (intra-cluster similarity) in its formulation. These two issues make (D)SE render indistinctive entropy to detect hallucination.

## 4 Incorporating intra-and inter-cluster similarity in uncertainty quantification

We can mitigate the information missing issues of (D)SE by leveraging both the intra-and inter-cluster similarities between generated answers, leading to more accurate uncertainty estimation. To alleviate the effect of outliers, we leverage LogSumExp operation to aggregate similarity. Putting together, *without the need of clustering*, we define Semantic Nearest Neighbor Entropy (SNNE) as

$$\text{SNNE}(\boldsymbol{q}) = -\frac{1}{n} \sum_{i=1}^n \log \sum_{j=1}^n \exp\left(\frac{f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q})}{\tau}\right) \qquad (3)$$

where $f$ measures the similarity between two answers given the question. The inner summation in Eq 3 effectively accounts for both intra- and inter-cluster similarities without requiring clustering, as SE does. Instead of explicitly grouping outputs into clusters, $f$ naturally captures intra-cluster similarity when $\boldsymbol{a}^i$ and $\boldsymbol{a}^j$ belong to the same semantic group and inter-cluster similarity when they do not. Because LogSumExp operation is a smooth approximation to the maximum function, SNNE resembles the entropy estimation based on the nearest neighbor distances (Beirlant et al., 1997). Thus, our method is less sensitive to outliers compared to LexSim and graph-based approaches.

**Extending SNNE to white-box settings.** Token probabilities of the generated outputs can be incorporated to weight the summation in Eq 3, enhancing the method's sensitivity to model confidence. Based on this, we propose the white-box version of SNNE as:

$$\text{WSNNE}(\boldsymbol{q})$$

$$= -\sum_{i=1}^{n} \bar{P}(\boldsymbol{a}^i|\boldsymbol{q}) \log \sum_{j=1}^{n} \exp\left(\frac{f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q})}{\tau}\right) \tag{4}$$

where $\bar{P}(\boldsymbol{a}^i|\boldsymbol{q}) = \tilde{P}(\boldsymbol{a}^i|\boldsymbol{q})/\sum_{j=1}^{n} \tilde{P}(\boldsymbol{a}^j|\boldsymbol{q})$

**Choices of similarity function.** The similarity function $f$ can be selected from several options. The first choice is ROUGE-L (Lin, 2004) used in LexSim (Fomicheva et al., 2020). Another option is to use the predicted scores from NLI models, as in graph-based methods (Lin et al., 2024). Finally, $f$ can also be defined as the cosine similarity between sentence embeddings of the model outputs.

**Theoretical results.** The following theorem shows that if we neglect inter-cluster similarity by assigning zero similarity to answers from different semantic clusters and intra-cluster similarity by assigning a constant value to answers within the same clusters, SNNE recovers DSE.

**Theorem 4.1.** *Let function $f$ be defined as $f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q}) = \tau \log(\frac{1}{n})$ if $\boldsymbol{a}^i, \boldsymbol{a}^j \in C_k$, otherwise $-\infty$. SNNE is equivalent to DSE as defined in Eq 2.*

Instead of a constant value, if we assign a term based solely on the sequence probability, WSNNE recovers SE.

**Theorem 4.2.** *Let function $f$ be defined as $f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q}) = \tau \log(\frac{\tilde{P}(\boldsymbol{a}^j|\boldsymbol{q})}{Q})$ if $\boldsymbol{a}^i, \boldsymbol{a}^j \in C_k$, otherwise $-\infty$ where $Q = \sum_{i=1}^{n} \tilde{P}(\boldsymbol{a}^i|\boldsymbol{q})$. WSNNE is equivalent to SE as defined in Eq 1.*

The proofs of Theorems 4.1 and 4.2 are provided in Appendix A. As stated, (D)SE is a special case of (W)SNNE under specific similarity metrics. Specifically, Theorem 4.1 shows that when a constant similarity is assigned to examples within a cluster and inter-cluster interactions are ignored, SNNE simplifies to DSE. Similarly, Theorem 4.2 demonstrates that when intra-cluster interactions are disregarded, WSNNE reduces to SE. These results imply that (W)SNNE provides a more expressive measure of uncertainty than (D)SE, and that leveraging informative similarity metrics enables (W)SNNE to outperform (D)SE, as evidenced by our experiments. A similar form of generalization was also observed in KLE's framework (Nikitin et al., 2024).

Furthermore, our method can asymptotically ($\tau \to 0$) recover LUQ-Pair when $f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q}) = \tau P(\text{entail}|\boldsymbol{a}^i, \boldsymbol{a}^j)$. Therefore, our method can be integrated into their method to find better atomic scores. Empirically, we showed that SNNE outperforms LUQ-PAIR for one-sentence generation, highlighting the better atomic uncertainty estimation of our method.

# 5 Experiments

## 5.1 Settings

**Models.** We use Llama-3.1-8B (Dubey et al., 2024) and Phi-3-mini-4k-instruct (Abdin et al., 2024).

**Datasets.** We evaluate our method across three NLP tasks: question answering (QA), text summarization (TS), and machine translation (MT). For QA, we use five datasets following SE (Farquhar et al., 2024) while for TS and MT, we adopt two summarization and two translation datasets in LM-Polygraph (Fadeeva et al., 2023). Details are deferred to Appendix B.

**Evaluation metrics.** For QA tasks, we assess UQ methods using AUROC and AUARC (Nadeem et al., 2009), following Lin et al. (2024). For TS and MT tasks, we adopt PRR (Malinin et al., 2017) in line with LM-Polygraph.

**Baselines.** We compare our methods with white-box and black-box UQ methods in Section 2.

Additional details are given in Appendix B.

## 5.2 Results

**QA datasets.** Figure 2 shows the performance of various uncertainty quantification (UQ) methods on answers generated by Llama-3.1-8B and Phi-3-mini. Both SNNE and WSNNE consistently outperform existing white-box and black-box baselines by a significant margin across both models. Among white-box methods, KLE$_{\text{full}}$ ranks second on Llama-3.1-8B, yet its black-box counterpart, KLE$_{\text{heat}}$, performs worse than the simpler DSE. Notably, our approach surpasses SAR, the previous state-of-the-art on the LM-Polygraph benchmark (Vashurin et al., 2025), especially in cases involving longer, single-sentence outputs where SAR's performance declines.

**TS and MT datasets.** Figure 3 demonstrates the superiority of our methods in terms of PRR score. LexSim is the runner-up in the black-box setting though its QA performance is much worse than the other baselines. For very long generations, the number of overlapping words in LexSim is a good indicator of semantic similarity as we confirmed in Table 5 in Appendix C. Notably, our methods improve over (D)SE and KLE by a clear margin.
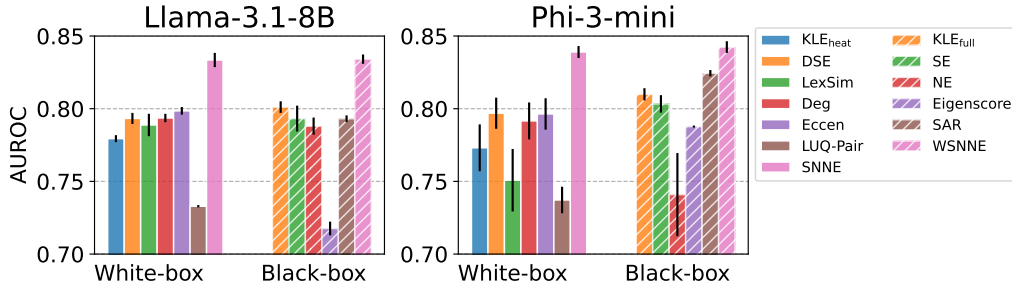
Figure 2: Average AUROC score of Llama-3.1-8B and Phi-3-mini on 5 QA tasks.
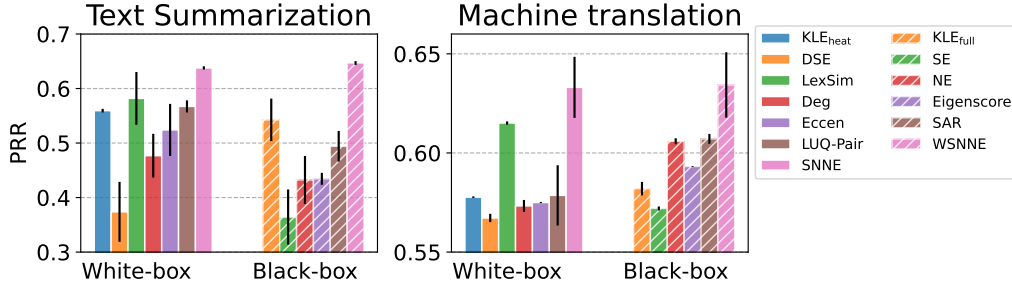


Figure 3: Average PRR score of Phi-3-mini on 2 text summarization and 2 machine translation tasks.
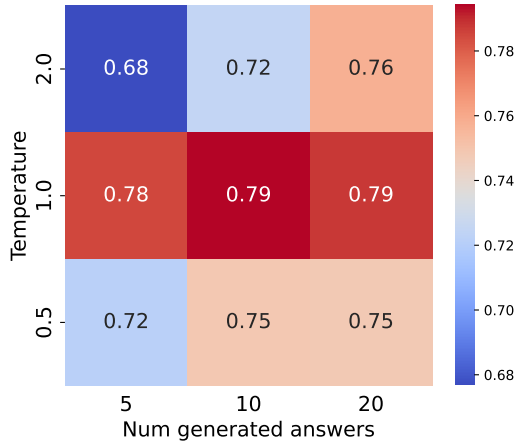


Figure 4: Effect of number of generated answers and generation temperature on the performance of SNNE. We measure the AUROC score of Llama-3.1-8B on SQuAD dataset.

This highlights the importance of incorporating intra-and inter-cluster similarities explicitly into entropy estimation in the long-generation cases.

**Effect of number of generations and generation temperature.** Figure 4 shows the impact of the number of model outputs and temperature on the AUROC score of SNNE. Increasing the number of generations enhances performance but comes at a higher inference cost. Additionally, at low temperatures, performance saturates at around 10 outputs. Using either too low (0.5) or too high (2.0) temperature degrades performance, as overly conservative or excessively diverse outputs compromise the quality of entropy estimation. Note

that our generation setting (temperature 1.0 and 10 generated answers) is directly adopted from SE's paper (Farquhar et al., 2024).

**Additional results.** Results for other models and ablation studies on the similarity function, scale factor $\tau$ can be found in Appendix C.

## 6 Conclusion

In this paper, we introduced a novel black-box uncertainty quantification method to address the limitations of semantic entropy in detecting hallucinations in long-generation scenarios. Our approach effectively accounts for both intra-and inter-cluster similarities and mitigate outliers, which are critical for accurate uncertainty estimation in modern LLMs. We also demonstrated that our method can be extended to white-box settings and provided theoretical results showing its generalization of semantic entropy. Extensive experiments on multiple LLMs and text generation tasks show that our method consistently outperforms existing uncertainty quantification methods.

## 7 Limitations

In this paper, we did not investigate uncertainty estimation in cases where the model generates multiple sentences or an entire paragraph. A naive approach would be to compute the uncertainty for each sentence independently and then aggregate these values into a single scalar. We leave this direction for our future work. Additionally, for differ-

ent data formats such as mathematical expressions, LaTeX equations, or code, our method requires further considerations. Designing an appropriate similarity function could help generalize our approach to these types of data. Finally, our method, similar to other existing UQ methods, requires sample multiple answers to estimate entropy, incurring additional inference cost.

## Acknowledgements

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Jan Beirlant, Edward J Dudewicz, László Györfi, Edward C Van der Meulen, et al. 1997. Nonparametric entropy estimation: An overview. *International Journal of Mathematical and Statistical Sciences*, 6(1):17–39.

Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, et al. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*, pages 12–58.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.

Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024. Inside: Llms' internal states retain the power of hallucination detection. *arXiv preprint arXiv:2402.03744*.

Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2023. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. *arXiv preprint arXiv:2307.01379*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2023. LM-polygraph: Uncertainty estimation for language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 446–461, Singapore. Association for Computational Linguistics.

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630.

Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Hsiu-Yuan Huang, Yutong Yang, Zhaoxi Zhang, Sanwoo Lee, and Yunfang Wu. 2024. A survey of uncertainty estimation in llms: Theory meets practice. *arXiv preprint arXiv:2410.15326*.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.

AQ Jiang, A Sablayrolles, A Mensch, C Bamford, DS Chaplot, D de las Casas, F Bressand, G Lengyel, G Lample, L Saulnier, et al. 2023. Mistral 7b (2023). *arXiv preprint arXiv:2310.06825*.

Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.

Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. Bioasq-qa: A manually curated corpus for biomedical question answering. *Scientific Data*, 10(1):170.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.

Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. Generating with confidence: Uncertainty quantification for black-box large language models. *Transactions on Machine Learning Research*.

Andrey Malinin and Mark Gales. 2020. Uncertainty estimation in autoregressive structured prediction. *arXiv preprint arXiv:2002.07650*.

Andrey Malinin, Anton Ragni, Kate Knill, and Mark Gales. 2017. Incorporating uncertainty into deep learning for spoken language assessment. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 45–50, Vancouver, Canada. Association for Computational Linguistics.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.

Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.

Malik Sajjad Ahmed Nadeem, Jean-Daniel Zucker, and Blaise Hanczar. 2009. Accuracy-rejection curves (arcs) for comparing classification methods with a reject option. In *Machine Learning in Systems Biology*, pages 65–81. PMLR.

Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*.

Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Fine-grained uncertainty quantification for llms from semantic similarities. *arXiv preprint arXiv:2405.20003*.

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*.

Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*.

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024a. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024b. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Daniil Vasilev, Akim Tsvigun, Sergey Petrakov, Rui Xing, Abdelrahman Sadallah, Kirill Grishchenkov, et al. 2025. Benchmarking uncertainty quantification methods for large language models with lm-polygraph. *Transactions of the Association for Computational Linguistics*, 13:220–248.

Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024. Luq: Long-text uncertainty quantification for llms. *arXiv preprint arXiv:2403.20279*.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International conference on machine learning*, pages 11328–11339. PMLR.

Rui Zhang and Joel Tetreault. 2019. This email could save your life: Introducing the task of email subject line generation. *arXiv preprint arXiv:1906.03497*.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

## A  Proofs

### A.1  Proof of Theorem 4.1

*Proof.* We define the similarity function $f$ as

$$f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q}) = \begin{cases} \tau \log(\frac{1}{n}), & \text{if } \exists k \text{ s.t. } \boldsymbol{a}^i, \boldsymbol{a}^j \in C_k \\ -\infty, & \text{otherwise} \end{cases} \tag{5}$$

In other words, we give a constant similarity for answers belonging to the same semantic group and zero similarity for answers in different groups. Then, the total similarity between $\boldsymbol{a}^i$ to all answers $\{\boldsymbol{a}^j\}_{j=1}^n$ becomes $\frac{|C_k|}{n}$ where $C_k$ is the semantic group of $\boldsymbol{a}^i$. Therefore, Eq 3 reads

$$\text{SNNE}(\boldsymbol{q}) = -\frac{1}{n} \sum_{i=1,\boldsymbol{a}^i \in C_k}^{n} \log \frac{|C_k|}{n} \tag{6}$$

$$= -\sum_{k=1}^{M} \frac{|C_k|}{n} \log \frac{|C_k|}{n} \tag{7}$$

$$= \text{DSE}(\boldsymbol{q}) \tag{8}$$

The penultimate equality holds because $n = \sum_{k=1}^{M} |C_k|$. $\qquad \square$

### A.2  Proof of Theorem 4.2

*Proof.* For a given question $\boldsymbol{q}$, let $Q$ be the total length-normalized sequence probability of its generated answers. We have the following equalities.

$$Q = \sum_{i=1}^{n} \tilde{P}(\boldsymbol{a}^i|\boldsymbol{q}) \tag{9}$$

$$= \sum_{k=1}^{M} P(C_k). \tag{10}$$

We define the similarity function $f$ as

$$f(\boldsymbol{a}^i, \boldsymbol{a}^j|\boldsymbol{q}) = \begin{cases} \tau \log(\frac{\tilde{P}(\boldsymbol{a}^j|\boldsymbol{q})}{Q}), & \text{if } \boldsymbol{a}^i, \boldsymbol{a}^j \in C_k \\ -\infty, & \text{otherwise} \end{cases} \tag{11}$$

Similar to Eq 5, we give a zero similarity for answers in different groups but give an asymmetric similarity based on the length-normalized sequence probability for answers belonging to the same semantic group. Then, the total similarity between $\boldsymbol{a}^i$ to all answers $\{\boldsymbol{a}^j\}_{j=1}^n$ becomes $\frac{P(C_k)}{Q} = \bar{P}(C_k)$ where $C_k$ is the semantic group of $\boldsymbol{a}^i$. Therefore,

Eq 4 reads

$$\text{SNNE}(\boldsymbol{q}) = -\sum_{i=1,\boldsymbol{a}^i \in C_k}^{n} \bar{P}(\boldsymbol{a}^i|\boldsymbol{q}) \log \bar{P}(C_k) \tag{12}$$

$$= -\sum_{i=1,\boldsymbol{a}^i \in C_k}^{n} \frac{\tilde{P}(\boldsymbol{a}^i|\boldsymbol{q})}{Q} \log \bar{P}(C_k) \tag{13}$$

$$= -\sum_{k=1}^{M} \frac{P(C_k)}{Q} \log \bar{P}(C_k) \tag{14}$$

$$= -\sum_{k=1}^{M} \bar{P}(C_k) \log \bar{P}(C_k) \tag{15}$$

$$= \text{SE}(\boldsymbol{q}) \tag{16}$$

$\qquad \square$

## B  Additional experimental details

**Models.** We use Llama2-7B, Llama2-13B (Touvron et al., 2023), Llama-3.1-8B (Dubey et al., 2024), Phi-3-mini-4k-instruct (Abdin et al., 2024), gemma-2-2b-it (Team et al., 2024b), and Mistral-Nemo-Instruct-2407 (Jiang et al., 2023).

**Datasets.** We evaluate our method across three NLP tasks: question answering (QA), text summarization (TS), and machine translation (MT). Following SE (Farquhar et al., 2024), we cover four different QA categories: commonsense knowledge (SQuAD (Rajpurkar et al., 2018) and TriviaQA (Joshi et al., 2017)), general knowledge from Google search (NaturalQuestion (Kwiatkowski et al., 2019)), simple math problems (Svamp (Patel et al., 2021)), biology and medicine (BioASQ (Krithara et al., 2023)). For text summarization, we adopt the abstractive single-document summarization dataset (XSUM (Narayan et al., 2018)) and the email subject line generation dataset (AESLC (Zhang and Tetreault, 2019)), as in LM-Polygraph (Fadeeva et al., 2023). For machine translation, we evaluate on two widely used datasets: WMT-14 German-to-English and WMT-14 French-to-English (Bojar et al., 2014).

**Baselines.** We compare our methods with white-box UQ methods including $\text{KLE}_{\text{full}}$ (Nikitin et al., 2024), SE (Farquhar et al., 2024), NE (Malinin and Gales, 2020), Eigenscore (Chen et al., 2024), SAR (Duan et al., 2023) and black-box UQ methods including $\text{KLE}_{\text{heat}}$ (Nikitin et al.,

Table 1: Instruction prompts for different tasks.

| Task | Instruction |
|---|---|
| QA | Answer the following question as briefly as possible. |
| XSUM | Here's the text and it's short one-sentence summary. |
| AESLC | Write a short subject line for the email. Output only the subject line itself. |
| WMT-14 de-en | Here is a sentence in German language and its translation in English language. |
| WMT-14 fr-en | Here is a sentence in French language and its translation in English language. |

Table 2: Average AUROC score on 5 QA tasks.

| Model | Acc | White-box | | | | Black-box | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $KLE_{full}$ | SE | NE | WSNNE | $KLE_{heat}$ | DSE | pTrue | NumSet | LexSim | SumEigv | Deg | Eccen | SNNE |
| Llama2-7B | 41.73 | 0.80 | 0.79 | 0.73 | **0.81** | 0.78 | 0.79 | 0.67 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | **0.80** |
| Llama2-13B | 46.07 | 0.79 | 0.78 | 0.74 | **0.80** | 0.78 | 0.78 | 0.74 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | **0.80** |
| Llama-3.1-8B | 50.76 | 0.80 | 0.79 | 0.79 | **0.83** | 0.78 | 0.79 | 0.72 | 0.79 | 0.79 | 0.79 | 0.79 | 0.80 | **0.83** |
| Phi-3-mini | 42.84 | 0.81 | 0.80 | 0.74 | **0.84** | 0.77 | 0.80 | 0.56 | 0.79 | 0.75 | 0.77 | 0.79 | 0.80 | **0.84** |
| gemma-2-2b | 39.28 | 0.83 | 0.82 | 0.77 | **0.84** | 0.82 | 0.82 | 0.72 | 0.82 | 0.80 | 0.81 | 0.82 | 0.82 | **0.84** |
| Mistral-Nemo | 54.02 | 0.80 | 0.79 | 0.75 | **0.82** | 0.78 | 0.79 | 0.74 | 0.79 | 0.78 | 0.78 | 0.79 | 0.79 | **0.81** |

Table 3: Average AUARC score on 5 QA tasks.

| Model | Acc | White-box | | | | Black-box | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $KLE_{full}$ | SE | NE | WSNNE | $KLE_{heat}$ | DSE | pTrue | NumSet | LexSim | SumEigv | Deg | Eccen | SNNE |
| Llama2-7B | 41.73 | **0.63** | 0.61 | 0.59 | **0.63** | 0.62 | 0.61 | 0.54 | 0.61 | 0.61 | 0.62 | 0.61 | 0.62 | **0.63** |
| Llama2-13B | 46.07 | **0.67** | 0.65 | 0.63 | 0.66 | 0.65 | 0.64 | 0.63 | 0.63 | 0.66 | 0.64 | 0.63 | 0.65 | **0.67** |
| Llama-3.1-8B | 50.76 | 0.73 | 0.71 | 0.72 | **0.74** | 0.71 | 0.71 | 0.66 | 0.70 | 0.71 | 0.71 | 0.71 | 0.71 | **0.74** |
| Phi-3-mini | 42.84 | 0.64 | 0.63 | 0.60 | **0.66** | 0.63 | 0.63 | 0.45 | 0.63 | 0.61 | 0.62 | 0.63 | 0.63 | **0.66** |
| gemma-2-2b | 39.28 | **0.64** | 0.62 | 0.59 | 0.63 | 0.62 | **0.63** | 0.54 | 0.62 | 0.60 | 0.61 | 0.62 | 0.62 | **0.63** |
| Mistral-Nemo | 54.02 | **0.75** | 0.72 | 0.70 | 0.74 | 0.72 | 0.72 | 0.70 | 0.72 | 0.72 | 0.72 | 0.72 | **0.73** | **0.73** |

Table 4: Average PRR score on 2 summarization tasks (XSUM, AESLC) and 2 translation tasks (WMT-14 de-en, WMT-14 fr-en). For correctness metric, R denotes ROUGE-L and B denotes BERTScore.

| Model | Metric | Score | White-box | | | | Black-box | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $KLE_{full}$ | SE | NE | WSNNE | $KLE_{heat}$ | DSE | pTrue | NumSet | LexSim | SumEigv | Deg | Eccen | SNNE |
| | | | | | | | Summarization | | | | | | | | |
| Phi-3-mini | R | 0.11 | 0.20 | 0.14 | 0.17 | **0.27** | 0.21 | 0.14 | 0.24 | 0.17 | 0.23 | 0.19 | 0.19 | 0.20 | **0.26** |
| | B | 0.44 | 0.54 | 0.36 | 0.43 | **0.65** | 0.56 | 0.37 | 0.52 | 0.38 | 0.58 | 0.47 | 0.48 | 0.52 | **0.64** |
| | | | | | | | Translation | | | | | | | | |
| Phi-3-mini | R | 0.62 | 0.58 | 0.57 | 0.60 | **0.63** | 0.58 | 0.57 | 0.56 | 0.56 | 0.61 | 0.57 | 0.57 | 0.58 | **0.63** |
| | B | 0.92 | 0.71 | 0.70 | 0.73 | **0.75** | 0.71 | 0.70 | 0.70 | 0.70 | 0.73 | 0.70 | 0.71 | 0.71 | **0.75** |

Table 5: Effect of similarity function of SNNE 2 summarization tasks (XSUM, AESLC) and 2 translation tasks (WMT-14 de-en, WMT-14 fr-en).

| Model | Metric | Translation | | | Summarization | | |
|---|---|---|---|---|---|---|---|
| | | embed | entail | ROUGE-L | embed | entail | ROUGE-L |
| Phi-3-mini | ROUGE-L | $0.62 \pm 0.02$ | $0.60 \pm 0.02$ | $\mathbf{0.63} \pm 0.02$ | $0.22 \pm 0.01$ | $0.22 \pm 0.01$ | $\mathbf{0.26} \pm 0.02$ |
| | BERTScore | $\mathbf{0.75} \pm 0.01$ | $0.73 \pm 0.01$ | $\mathbf{0.75} \pm 0.01$ | $0.56 \pm 0.02$ | $0.57 \pm 0.01$ | $\mathbf{0.64} \pm 0.01$ |

2024), DSE (Farquhar et al., 2024), pTrue (Kadavath et al., 2022), LexSim (Fomicheva et al., 2020), graph-based methods (NumSet, SumEigv, Deg, Eccen) (Lin et al., 2024), LUQ-Pair (Zhang et al., 2024).

**Generation setting.** For each question, we generate one answer at a low temperature ($T = 0.1$) to assess model correctness and sample 10 answers at a high temperature ($T = 10$) to estimate uncertainty. We use 5-shot in-context demonstrations for QA tasks and instruct the model to produce short answers following the brief prompt from SE. For TS and MT tasks, we employ a 0-shot setting and adopt the instructions from LM-Polygraph. Table 1 provides a summary of the prompts used across different tasks.

**Correctness measures.** Following (Farquhar et al., 2024), we evaluate QA tasks using the F1

| Model | Task | Metric | $\tau$ | WSNNE | SNEE |
|---|---|---|---|---|---|
| Llama-3.1-8B | QA | AUROC | 0.1 | $0.830 \pm 0.003$ | $0.830 \pm 0.004$ |
| | | | 1 | $\mathbf{0.833 \pm 0.002}$ | $\mathbf{0.832 \pm 0.004}$ |
| | | | 10 | $0.832 \pm 0.003$ | $0.830 \pm 0.004$ |
| | | | 100 | $0.831 \pm 0.003$ | $0.830 \pm 0.004$ |
| Phi-3-mini | QA | AUROC | 0.1 | $0.836 \pm 0.003$ | $0.832 \pm 0.003$ |
| | | | 1 | $\mathbf{0.841 \pm 0.004}$ | $\mathbf{0.838 \pm 0.004}$ |
| | | | 10 | $0.840 \pm 0.004$ | $0.836 \pm 0.004$ |
| | | | 100 | $0.840 \pm 0.004$ | $0.834 \pm 0.004$ |
| Phi-3-mini | TS | PRR | 0.1 | $0.645 \pm 0.002$ | $0.634 \pm 0.002$ |
| | | | 1 | $\mathbf{0.646 \pm 0.005}$ | $\mathbf{0.636 \pm 0.005}$ |
| | | | 10 | $0.645 \pm 0.005$ | $0.635 \pm 0.005$ |
| | | | 100 | $0.645 \pm 0.005$ | $0.635 \pm 0.005$ |
| Phi-3-mini | MT | PRR | 0.1 | $0.632 \pm 0.015$ | $0.631 \pm 0.015$ |
| | | | 1 | $\mathbf{0.635 \pm 0.016}$ | $\mathbf{0.632 \pm 0.016}$ |
| | | | 10 | $0.634 \pm 0.016$ | $0.632 \pm 0.016$ |
| | | | 100 | $0.634 \pm 0.016$ | $0.632 \pm 0.016$ |

Table 6: Effect of the scale factor $\tau$ on 2 summarization tasks ( XSUM , AESLC ) and 2 translation tasks (WMT-14 de-en, WMT-14 fr-en).

score, a standard metric for the SQuAD dataset, with a correctness threshold of 50% applied across all QA datasets. For text summarization and machine translation, we use ROUGE-L (Lin, 2004) and BERTScore (Zhang* et al., 2020), as in LM-Polygraph.

**Hyperparameter setting.** We select the best scale factor $\tau$ in Eq 3 and 4 from $\{0.1, 1, 10, 100\}$. For similarity function $f$, we use ROUGE-L score.

**Implementation details.** For calculating ROUGE-L, we use the Python implemention of Google https://github.com/google-research/google-research/tree/master/rouge. For models and datasets, we download them from Hugging Face https://huggingface.co/.

**Computational resources.** We conduct each experiment three times using NVIDIA RTX A6000 GPUs.

## C  Additional experimental results

### C.1  Detailed results

**QA datasets.** Table 2 demonstrates the average AUROC score of 6 different models on 5 QA datasets. Across different model architectures, SNNE and WSNNE are the best-performing white-box and black-box methods, respectively. KLE$_{\text{full}}$ is the second-best white-box method but its black-box KLE$_{\text{heat}}$ is surpassed by DSE which is

the runner-up in the black-box setting. In terms of the AUARC score, Table 3 shows that SNNE still yields the highest score in the black-box setting while KLE$_{\text{full}}$ exhibits a competitive performance with WSNNE.

**TS and MT datasets.** Table 4 illustrates the average PRR score with respect to two different correctness measure: ROUGE-L and BERTScore. In the black-box scenario, SNNE shows the superior advantages over other baselines across different settings, followed by LexSim. By capturing the semantic distance in the kernel space, KLE improves over (D)SE in both white-box and black-box settings though they are still behind our methods by a clear margin. This reinforces the effectiveness of explicitly integrating the intra-and inter-cluster semantic similarities into entropy calculation.

### C.2  Ablation studies

**Effect of similarity function on SNNE.** We examine three different choices of similarity functions: ROUGE-L, entail, and embed in Section 4. For entail with NLI models, we use the same DeBERTa model (He et al., 2020) that is used for semantic clustering. For embed with sentence transformer, we use the Qwen2-7B-instruct model (Li et al., 2023) which ranked 1st in the MTEB benchmark (Muennighoff et al., 2022). Ta-

ble 5 summarizes the PRR score on TS and MT tasks. ROUGE-L results in the best performance except for translation task with BERTScore in which ROUGE-L ties with embed. The superior of ROUGE-L also explains the good performance of LexSim in Table 4.

**Effect of scale factor $\tau$.** When $\tau$ is large, SNNE smooths out differences, making all intra-and inter-distances contribute more equally. In that case, SNNE behaves more like an average pairwise similarity similar to Deg or LexSim. In contrast, when $\tau \to 0$, it emphasizes the smallest "intra-distance" and ignores all other intra-and inter-distances. Additionally, in one-sentence output scenarios, and with a specific design of the similarity function $f$, our method asymptotically ($\tau \to 0$) recovers LUQ-Pair. Table 6 illustrates the impact of varying scale factor on QA tasks in Figure 2 and TS & MT in Figure 3 below. Our method is not sensitive to the choice of $\tau$ in the tuning set. Overall, $\tau = 1$ yields the best performance across different use cases, making it our default choice.