

API Is Enough: Conformal Prediction for Large Language Models Without Logit-Access

Anonymous ACL submission

Abstract

This study aims to address the pervasive challenge of quantifying uncertainty in large language models (LLMs) without logit-access. Conformal Prediction (CP), known for its model-agnostic and distribution-free features, is a desired approach for various LLMs and data distributions. However, existing CP methods for LLMs typically assume access to the logits, which are unavailable for some API-only LLMs. In addition, logits are known to be miscalibrated, potentially leading to degraded CP performance. To tackle these challenges, we introduce a novel CP method that (1) is tailored for API-only LLMs without logit-access; (2) minimizes the size of prediction sets; and (3) ensures a statistical guarantee of the user-defined coverage. The core idea of this approach is to formulate nonconformity measures using both coarse-grained (i.e., sample frequency) and fine-grained uncertainty notions (e.g., semantic similarity). Experimental results on both close-ended and open-ended Question Answering tasks show our approach can mostly outperform the logit-based CP baselines.

1 Introduction

Large Language Models (LLMs) have made significant advancements (Thoppilan et al., 2022; Wei et al., 2022, 2023), highlighting the research potential of natural language generation (Peinl and Wirth, 2023). However, they often generate information that is not accurate, factual, or grounded in reality, referred to as "hallucination" (LeCun, 2023). Therefore, it is crucial to quantify LLM uncertainty to ensure responsible responses.

However, uncertainty quantification (UQ) for LLMs is challenging due to the complex data distributions and inner model mechanism, as well as the often limited access to logit information. A potential solution is to use conformal prediction (CP) (Vovk et al., 2005; Angelopoulos and Bates, 2021; Kato et al., 2023; Wang et al., 2023), which

is known for being model-agnostic and distribution-free, and with rigorous coverage guarantees. Given a user-defined error rate α , CP provides a guaranteed coverage rate for prediction sets/intervals. It measures the uncertainty from a model prediction using nonconformity score functions, e.g., $1 - f(X)_Y$ (Sadinle et al., 2019), where $f(X)_Y$ is the softmax score for the true label Y .

Most of the existing CPs for LLMs rely on the access to model logits to measure nonconformity scores. For instance, Kumar et al. (2023) define nonconformity scores as softmax scores for logits of different options in the multi-choice question answering (MCQ) task and Quach et al. (2023) apply the conformal risk control framework (Angelopoulos et al., 2021), an extension of CP, to LLMs by utilizing model-based log probability. However, for some API-only LLMs like Bard (Manyika and Hsiao, 2023), logit-access is almost impossible for end users. Even though the logits are available (e.g., for 4w (OpenAI, 2023)), they are known to be miscalibrated and can lead to degraded performance of CP w.r.t. estimating the prediction sets or intervals (Nguyen and O'Connor, 2015; Lin et al., 2022), e.g., a large set size (i.e., low efficiency).

To enable CP without logit-access, a straightforward way is to calculate the frequency of each response via sampling and approximate model-based probabilities. However, we theoretically prove that this approach is extremely computationally expensive (Lemma 3.1). As nonconformity scores typically measure the level of uncertainty, CP depends on the *ranking* of the nonconformity measures rather than their actual *values* (Shafer and Vovk, 2008). Therefore, we propose to sample responses for a certain number of times (e.g., 30) for each input and then utilize the frequency of each response as a *coarse-grained* uncertainty notion. This approach reduces the overall sampling costs and eliminates the dependence on the logits. However, when using frequency as the only non-

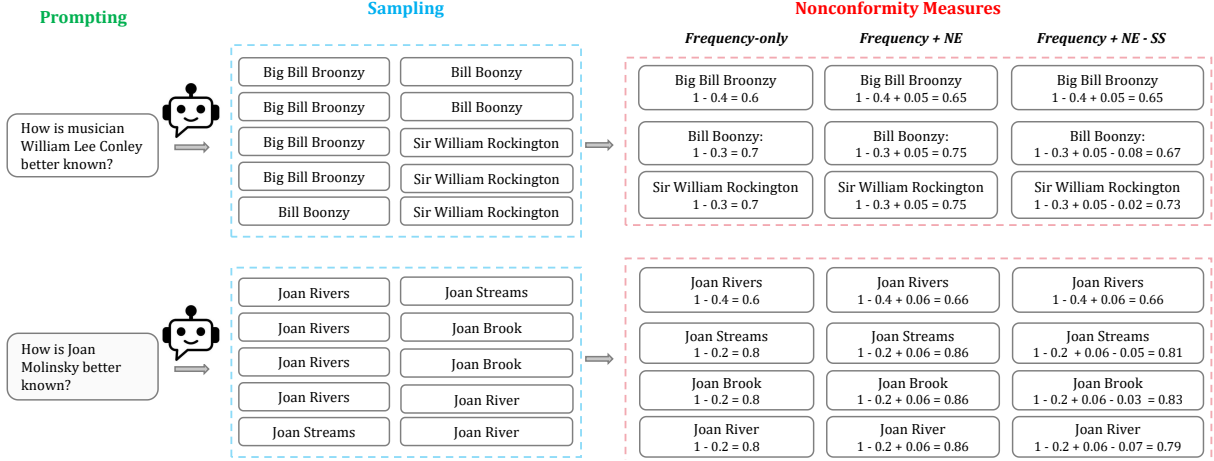


Figure 1: Illustrations of the proposed problem and solution. Three uncertainty notions for measuring nonconformity: (1) Frequency-only, where the nonconformity score is calculated as $1 - \text{the frequency of a response out of 10 samplings}$. Concentration issues arise at scores of **0.6, 0.7, and 0.8**. For instance, responses from different prompts (e.g., "Big Bill Broonzy" and "Joan Rivers") have the same score of 0.6, as well as responses within the same prompt (e.g., "Bill Boonzy" and "Sir William Rockington") which both have a score of 0.7, and so forth. (2) Frequency combined with NE, where the nonconformity score is calculated as $1 - \text{frequency} + \text{NE}$, revealing concentration issues at scores of **0.75 and 0.86**. (3) Frequency, NE, and SS combined, where the nonconformity score is calculated as $1 - \text{frequency} + \text{NE} - \text{SS}$, with **no observed concentration issues**.

conformity measure, we observe that nonconformity scores concentrate on certain values as some responses may share the same frequency even if they have varied levels of uncertainty (see Figure 1), consequently diminishing the efficiency of prediction sets.

To distinguish between responses that share the same frequency, we first identify two potential causes: the respective concentration issues across different prompts and within the same prompt, which indicates we need to integrate prompt-wise and response-wise notions to respectively mitigate these two causes. We then propose two additional *fine-grained* uncertainty notions: normalized entropy (NE), measuring prompt-wise self-consistency to alleviate concentration issues across different prompts; and semantic similarity (SS), measuring response-wise similarity to the most frequent response within the same prompt, to mitigate internal concentration issues specific to the prompt. Figure 1 illustrates the different scores defined using *frequency-only*, *frequency combined with NE*, and *frequency combined with NE and SS* as nonconformity measures, respectively. By considering various uncertainty information, the proposed nonconformity score function can better distinguish the uncertainty of different responses.

Our contributions are summarized as follows:

- To our knowledge, this is the first CP work dedicated to LLMs without logit-access that provides a coverage guarantee for the prediction set with a small size.

- We propose a novel CP approach that uses both course-grained and fine-grained uncertainty notions as the non-conformity measures. We also theoretically prove (1) it is computationally infeasible to use response frequency to approximate model output probability, and (2) our approach ensures a rigorous statistical coverage guarantee.
- We conduct experiments on both close- and open-ended QA tasks and demonstrate the effectiveness of our method. Notably, we mostly surpass all baselines, including four logit-access methods and one method without logit-access.

2 Preliminaries of Conformal Prediction

Conformal prediction (CP) (Vovk et al., 2005) is a model-agnostic method offering distribution-free uncertainty quantification, which produces prediction sets/intervals containing ground-truth labels with a desired error rate α . One of the widely used CP methods is split CP. Formally, let (X, Y) be a sample, where X represents features and Y represents the outcome. We denote the calibration set as $(X_i, Y_i)_{i=1, \dots, n}$ and the test set as $(X_{\text{test}}, Y_{\text{test}})$. CP presents the following nesting property:

$$\alpha_1 > \alpha_2 \Rightarrow C_{1-\alpha_1}(X) \subseteq C_{1-\alpha_2}(X). \quad (1)$$

where $C : X \rightarrow 2^Y$ is a set-valued function that generates prediction sets over the powerset of Y given an input X .

Theorem 2.1 (Conformal coverage guarantee). *Suppose $(X_i, Y_i)_{i=1, \dots, n}$ and $(X_{\text{test}}, Y_{\text{test}})$ are independent and identically distributed (i.i.d.).*

$C_{1-\alpha}(X_{test})$ is a set-valued mapping satisfying the nesting property in Eq. 1. Then the following holds:

$$P(Y_{test} \in C_{1-\alpha}(X_{test})) \geq 1 - \alpha, \quad (2)$$

where $\alpha \in (0, 1)$ is the user-defined error rate.

Nonconformity Measures. The nonconformity measure N is a core element in CP. It measures uncertainty in the model’s output by assessing the deviation of a specific instance or output from patterns observed in the training data. Typically, we have logit access to models to measure nonconformity, e.g., $1 - f(X)_Y$. For LLMs, N is typically derived from the post-hoc logits.

Split CP Steps. Split CP typically involves four steps (Angelopoulos and Bates, 2021):

1. Establish heuristic uncertainty notions.
2. Define the nonconformity measures/score function $N(x, y) \in \mathbb{R}$.
3. Compute \hat{q} as the $\frac{[(n+1)(1-\alpha)]}{n}$ quantile of the nonconformity scores.
4. Use \hat{q} to generate prediction sets for new examples: $C(X_{test}) = \{Y : N(X_{test}, Y) \leq \hat{q}\}$.

3 Methodology

Our method considers two pivotal challenges arising from the LLMs without logit-access: how to approximate the logit information of LLMs; and how to further improve CP efficiency, i.e., small prediction sets. We propose the **Logit-free Conformal Prediction for LLMs (LofreeCP)**, where its nonconformity measures consist of three notions: *frequency*, representing coarse-grained uncertainty; *NE*, representing prompt-wise fine-grained uncertainty; and *SS*, representing response-wise fine-grained uncertainty.

3.1 Frequency As the Rankings Proxy

A straightforward way is to approximate real predictive probabilities through a sufficiently large number of samplings. However, as we show in Lemma 3.1, a minimum of 9,604 samples is required to achieve a 95% confidence level with a 1% margin of error. Therefore, the implementation is impractical due to computational constraints.

Lemma 3.1 (Minimum Sample Size for Confident Probability Estimation). *Let $freq(Y_i)$ be the absolute frequency of outcome Y_i in the sampling, N_{total}*

be the total number of samplings, p_i be the desired estimated probability, ϵ be the estimation error, and δ be the target confidence level. To determine the minimum sample size for confident probability estimation, for any given $\epsilon > 0$ and $0 < \delta < 1$, the following inequality must hold:

$$P \left\{ \left| \frac{freq(Y_i)}{N_{total}} - p_i \right| \leq \epsilon \right\} \geq \delta. \quad (3)$$

Then, the minimum sample size N_{total} satisfying Inequality 3 is given by:

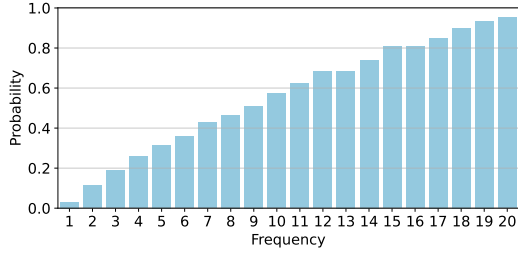
$$N_{total} \geq \left(\frac{u_{1-(1-\delta)/2}}{2\epsilon} \right)^2, \quad (4)$$

where $u_{1-(1-\delta)/2}$ is the quantile of the standard normal distribution corresponding to the confidence level $1 - (1 - \delta)/2$. The proof of Lemma 3.1 is given in Appendix A.1.

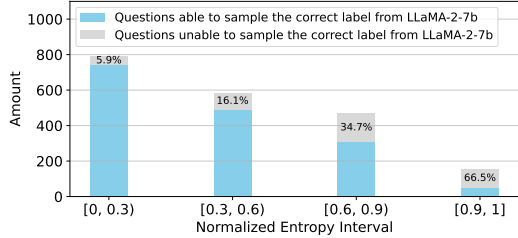
Since nonconformity measures are grounded in assessing the model’s predictive uncertainty (Shafer and Vovk, 2008), the primary focus lies in the *rankings* of uncertainty inherent in nonconformity measures rather than the absolute values themselves. Further, self-consistency theory (Wang et al., 2022; Li et al., 2022) states that a repetitively sampled response is viewed as a form of consistency linked to higher confidence in the response. To empirically validate this intuition, we randomly select 2000 questions from the TriviaQA dataset (Joshi et al., 2017). We conducted 20 samplings from the Llama-2-7b model (Touvron et al., 2023), extracted logits, and subsequently computed model output probabilities. The observed results depicted in Figure 2a indicate a direct positive correlation between response frequency and average real probability. As the response frequency climbs, there is a corresponding increase in the average real probability, suggesting a growing level of confidence and certainty in the model’s responses. Therefore, we propose to use frequency as the proxy of probability ranking. It is defined as

$$F(\hat{y}_a^{(i)}, m) = \frac{\tilde{p}[\hat{y}_a^{(i)}]}{m}, \quad (5)$$

where \tilde{p} represents the empirical absolute frequency, $\hat{y}_a^{(i)}$ is the a -th non-repeated sampled response for i -th prompt, m is the sampling quantity from LLMs for each prompt. However, only using response frequency as nonconformity measures results in the concentration of nonconformity scores on certain values. This issue makes it challenging to discern nonconformity differences among responses with the same scores, rendering ineffective calibration in CP.



(a) LLM response frequency vs. LLM output probability.



(b) NE vs. the proportion of prompts unable to sample correct labels from Llama-2-7b.

Figure 2: Empirical findings with TriviaQA dataset.

3.2 Fine-grained Uncertainty Notions

To resolve the concentration issue, we propose two fine-grained uncertainty measures. Firstly, inspired by self-consistency theory (Wang et al., 2022; Li et al., 2022), we incorporate NE, a prompt-wise fine-grained uncertainty notion, to mitigate the concentration issue across different prompts. NE is a measure of the uncertainty or diversity in the model’s predictions when generating responses to a given prompt. It is defined as

$$H(x^{(i)}|\{\hat{y}_j^{(i)}\}_{j=1}^m) = \left| \frac{\sum_{a=1}^n \tilde{F}(\hat{y}_a^{(i)}) \log(\tilde{F}(\hat{y}_a^{(i)}))}{\log m} \right|, \quad (6)$$

where $x^{(i)}$ is the i -th instance of the prompt dataset, m is the number of sampled responses, n is the number of non-repeated responses, $\hat{y}_j^{(i)}$ is the j -th sampled response. Following experiments in Section 3.1, we show that as NE increases, the number of unanswered questions also increases (Figure 2b), indicating a rise in uncertainty.

Secondly, to address concentration issues within a prompt, we introduce SS as a response-wise fine-grained uncertainty measure. This metric semantically assesses the similarity between each non-top-1 response and the top-1 response within a prompt. Intuitively, when two non-top-1 responses share the same frequency, the one more semantically similar to the top-1 response is more likely to express high confidence and low uncertainty. We use the cosine

similarity to express SS. It is defined as

$$SS(\hat{y}_a^{(i)}, P_{\text{highest}}^{(i)}) = \frac{\mathbf{v}(\hat{y}_a^{(i)}) \cdot \mathbf{v}(P_{\text{highest}}^{(i)})}{\|\mathbf{v}(\hat{y}_a^{(i)})\| \cdot \|\mathbf{v}(P_{\text{highest}}^{(i)})\|}, \quad (7)$$

where $\mathbf{v}(x)$ is the vector representation of x , $P_{\text{highest}}^{(i)}$ is the response having the highest frequency for i -th prompt. However, if the response to be measured is the one with the highest frequency, we do not consider SS with itself.

3.3 CP for LLMs Without Logit-Access

Considering both the coarse-grained and fine-grained uncertainty notions, the final nonconformity score function of LofreeCP is defined as

$$N^{(i)} = -F(\hat{y}_a^{(i)}, m) + \lambda_1 \cdot H(x^{(i)}|\{\hat{y}_j^{(i)}\}_{j=1}^m) - \lambda_2 \cdot SS(\hat{y}_a^{(i)}, P_{\text{highest}}^{(i)}), \quad (8)$$

where $\lambda = (\lambda_1, \lambda_2)$ representing a hyperparameter configuration controls the balance between the coarse-grained and fine-grained uncertainty notions. LofreeCP has the coverage guarantee:

Proposition 3.2 (Coverage guarantee of LofreeCP). *Suppose $(X_i, Y_i)_{i=1, \dots, n}$ and $(X_{\text{test}}, Y_{\text{test}})$ are i.i.d. Let $C_{1-\alpha}(X_{\text{test}})$ be defined as in Step 3. Then we have the coverage guarantee:*

$$P \{Y_{\text{test}} \in C_{1-\alpha}(X_{\text{test}})\} \geq 1 - \alpha,$$

where $\alpha \in (0, 1)$ denotes the desired error rate. The proof of the coverage guarantee of LofreeCP is provided in Appendix A.2.

LofreeCP consists of three stages: calibration, validation, and testing. The calibration stage aims to find the quantile based on the desired error rate. We sample m responses from the LLM for each prompt and store them in a response pool. Then, we obtain the nonconformity scores of the true labels with the following rules: if the true label exists in the pool, we use the nonconformity measures from Equation 8 to calculate its nonconformity score; otherwise, we set the nonconformity score as ∞ to signify that it is nearly impossible for the LLM to generate the true response. After obtaining all nonconformity scores of the calibration set, we find the quantile based on the desired error rate. We use this quantile as a threshold value for both the validation and test stages.

We then use the validation set to choose the optima hyperparameter configuration $\lambda = (\lambda_1, \lambda_2)$. Subsequently, we conduct evaluations on the test

set using the chosen configuration. Both stages follow identical sampling steps to the calibration, traversing all responses and calculating the nonconformity scores. We preserve the responses whose nonconformity score is less than the threshold in our final prediction set. The pseudocode of the LofreeCP method is provided in Appendix B.10.

4 Experiments

4.1 Experimental Setup

Backbone LLMs and Evaluation Tasks. Since we need to compare LofreeCP with logit-based methods, from where logits can be retrieved directly, we consider different open-source LLMs, including Llama-2-7B, Llama-2-13B, WizardLM-v1.2(13b) (Xu et al., 2023) and Vicuna-v1.5(7b) (Chiang et al., 2023) models as our backbone models. Note that our method uses these LLMs as if they were API-only LLMs, i.e., it assumes no access to any internal information of LLMs. We use both open-ended Question-Answering (QA) and close-ended Multi-Choice Question-Answering (MCQ) tasks for evaluation.

Datasets. We use standard benchmarking datasets TriviaQA and MMLU (Hendrycks et al., 2020), following (Kumar et al., 2023) and (Quach et al., 2023). We also include the WebQuestions benchmark (Berant et al., 2013). For QA, we use the TriviaQA dataset, which consists of trivia questions spanning a wide range of topics such as history and science, and the WebQuestions dataset, which is focused on questions asked by users on a search engine. MMLU dataset, covering 57 subjects (e.g., mathematics, history), is used for MCQ. We focus on a subset of 16 subjects out of the total 57, as in Kumar et al. (2023).

Baselines. Baselines include methods without logit-access and those based on logit:

- **Top- K_{white} .** A logit-based non-CP method without coverage guarantee, which includes responses with the first K highest probabilities for each prompt in the prediction set.
- **Standard Split Conformal Prediction (SCP)** (Vovk et al., 2005). A logit-based CP method, which follows the steps shown in Section 2.
- **Sorted Adaptive Prediction Sets (SAPS)** (Huang et al., 2023). A logit-based CP method, which uses the highest probability and replaces other probabilities with some weighted values to mitigate the miscalibration issue.

- **Top- K_{black} .** A non-CP method without logit-access and coverage guarantee, which includes responses with the first K highest frequency for each prompt in the prediction set.
- **Conformal Language Modeling (CLM)** (Quach et al., 2023). The state-of-the-art logit-based CP method, which uses the general risk control framework. This baseline is only used in QA as it is not applied to MCQ.

Metrics. We use following metrics for evaluation (Angelopoulos and Bates, 2021):

- Empirical Coverage Rate (ECR) assesses whether the conformal procedure has the correct coverage with the theoretical guarantee.
- Size-Stratified Coverage (SSC) (Angelopoulos et al., 2020) assesses the worst coverage rate of each bin among different set sizes.
- Average Prediction Set Size (APSS) assesses the efficiency of CP. We expect the APSS of an efficient CP method to be small.

4.2 Results for QA

We perform QA using TriviaQA and WebQuestions datasets. The results for Llama-2-13b are reported in Tables 1-2, those for Llama-2-7b are shown in the sensitivity analysis of Section 4.5 and those for WizardLM-v1.2(13b) and Vicuna-7b-v1.5 can be found in Appendix D. In Table 1, the LofreeCP method excels on TriviaQA across all error rate settings, outperforming the second-best method, CLM, by 7.7% in terms of APSS at an error rate of 0.25. Regarding SSC, our LofreeCP method surpasses the second-best method, First- K_{white} , by 1.6%. In Table 2, our method demonstrates superior performance on WebQuestions in most settings. For instance, at an error rate of 0.45, our LofreeCP method outperforms the second-best method, CLM, by 11.6% in terms of APSS. Regarding SSC, we outperform the second-best method, SCP, by 4.3%. WizardLM-v1.2(13b) and Vicuna-7b-v1.5 exhibit similar trends to Llama-2-13b.

The smallest APSS indicates that our method can produce the most efficient prediction sets. The highest SSC indicates that our method is attentive to the conditional coverage rate, achieving well-calibrated uncertainty estimates within diverse size categories. The rationale behind the observed superior performance is that our nonconformity measure can capture the coarse-grained uncertainty of responses and effectively optimize nonconformity through fine-grained considerations, thereby miti-

Table 1: Results for TriviaQA using Llama-2-13b: Among all baselines, only *First-K_{white}* and *First-K_{black}* are non-CP-based, while the rest are CP-based methods. In the results, **bold** indicates that the method produces the best performance among all methods; **X** denotes that the method fails to produce the set with the desired error rate.

Methods	Logit-Access	Error Rate								
		0.2			0.25			0.3		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	82.1	76.6	3.39	76.1	72.9	1.90	X	X	X
CLM	✓	80.2	73.4	2.29	75.2	69.1	1.55	70.1	68.3	1.28
SCP	✓	80.3	75.7	2.25	75.1	70.0	1.59	70.3	74.5	1.21
SAPS	✓	80.0	77.9	2.74	75.1	64.2	1.80	70.0	49.4	1.55
First-K _{black}	X	80.1	76.8	2.70	76.4	72.2	1.90	X	X	X
LofreeCP (Ours)	X	80.1	79.0	2.19	75.3	74.5	1.43	70.3	76.7	1.08

Methods	Logit-Access	Error Rate								
		0.35			0.4			0.45		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	X	X	X	62.4	62.5	1.00	X	X	X
CLM	✓	65.0	69.3	0.96	60.1	72.7	0.81	55.2	83.3	0.70
SCP	✓	65.1	76.4	1.02	60.3	75.7	0.85	55.3	82.5	0.74
SAPS	✓	65.1	57.4	1.28	60.1	70.7	0.85	55.1	76.5	0.72
First-K _{black}	X	66.5	66.5	1.00	X	X	X	X	X	X
LofreeCP (Ours)	X	65.1	78.5	0.90	60.0	81.0	0.75	55.2	84.1	0.66

Table 2: Results for WebQuestions using Llama-2-13b.

Methods	Logit-Access	Error rate											
		0.35			0.4			0.45			0.5		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	66.4	57.5	6.18	61.6	58.1	3.81	57.5	55.0	2.91	50.6	49.0	1.97
CLM	✓	65.3	50.5	4.54	60.5	52.9	2.86	55.0	51.6	1.81	50.1	56.8	1.27
SCP	✓	65.1	46.7	4.61	61.6	49.3	3.01	55.2	55.8	2.02	50.2	57.8	1.39
SAPS	✓	65.2	46.2	5.19	60.6	56.2	3.39	55.5	37.7	2.40	50.8	21.7	1.86
First-K _{black}	X	65.1	54.9	6.20	60.0	55.3	3.78	56.9	54.4	2.91	53.7	52.4	1.97
LofreeCP (Ours)	X	65.1	61.1	5.33	60.0	60.0	2.68	55.1	60.1	1.60	50.3	59.9	1.06

gating the inherent miscalibration issue in LLMs.

4.3 Ablation Study

To demonstrate the impact of our fine-grained uncertainty notions (NE and SS) on mitigating the concentration issues, we conduct a series of ablation studies using the TriviaQA dataset with a sampling quantity of 20. We compare LofreeCP with its different variants: we remove one fine-grained notion at a time (Freq&SS, removing the NE notion; and Freq&NE, removing the SS notion), and finally remove both fine-grained notions (Freq-Only). We report APSS and ECR, the direct indicators of the concentration issue, in Figure 3.

Impact of Concentration Issue. As introduced in Section 3, the concentration issue occurs when the nonconformity score is concentrated on certain values. When we use the frequency-only variant (Freq-Only), this issue can be observed in all error rate settings, as shown in Figure 3: Freq-Only has the largest APSS and the most conservative ECR. Due to its coarse-grained uncertainty notion, Freq-Only tends to generate similar nonconformity scores clustered into several groups, making it hard to differentiate granular uncertainties to produce efficient prediction sets.

Full Method Mitigates Concentration Issue.

We further observe that the concentration issue is mitigated in all error rate settings by incorporating fine-grained notions (NE & SS). For example, at an error rate of 0.2, Freq-Only exhibits an APSS of nearly 6.5, while the full method LofreeCP has an APSS of 4.27, resulting in a drop of more than 23%. The method including only SS or NE also mitigates the concentration issue to some extent, while the full method performs the best in terms of APSS and ECR. The results suggest that NE and SS both have a significant impact on improving the efficiency of prediction sets by mitigating concentration issues of nonconformity scores.

4.4 Results for MCQ

In addition to open-ended tasks, e.g. QA, LofreeCP is also effective at close-ended tasks that can be converted into a generation pipeline, e.g. MCQ. We conduct MCQ experiments on the MMLU dataset using Llama-2-13b with a sampling quantity of 20. We present the results in Figure 4.¹ LofreeCP exhibits superior performance. When compared with SCP and SAPS across all 16 subjects, LofreeCP

¹We omit the results from top-K methods as they exhibit much larger APSS than other methods for MCQ.

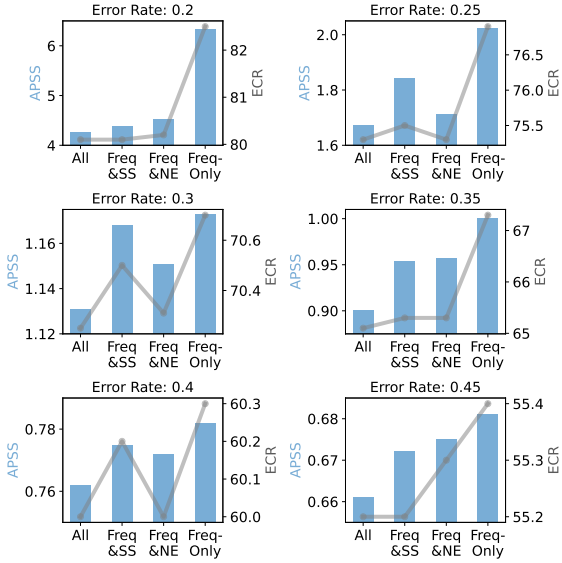


Figure 3: Ablation study. The blue bar chart represents APSS, while the gray line represents ECR.

achieves the best performance in 9 subjects and ties for the best in subjects of professional medicine, college chemistry, and marketing, resulting in the overall best performance in 12 out of 16 subjects. In contrast, SCP only ties for the best in 3 subjects. SAPS achieves the solo best performance in 3 subjects and ties for the best in 1 subject.

An intriguing observation is related to subjects in the business and management (B&M) category (e.g., marketing and public relations). When using LofreeCP method, these subjects show slightly larger APSS than the two logit-based methods, SCP and SAPS. This suggests that the logits for responses to B&M questions predicted by the Llama-2-13b model are better calibrated than the remaining subjects from the Science, Technology, Engineering, and Mathematics (STEM) category. Our LofreeCP method mitigates the model miscalibration issue by refraining from directly using logits.

4.5 Sensitivity Analyses

BackBone Models. To investigate the influence of different backbone models on the performance of LofreeCP, we conduct experiments using Llama-2-7b and Llama-2-13b with a sampling quantity of 20. Results of SSC and APSS are shown in Figure 5. We observe that better performance of APSS and SSC in the 13b setting than in the 7b setting. We believe this is because Llama-2-13b is more powerful than Llama-2-7b, and produces more confident and calibrated responses, thereby providing more efficient prediction sets. Results for Vicuna-v1.5(7b) are provided in Appendix D, indicating

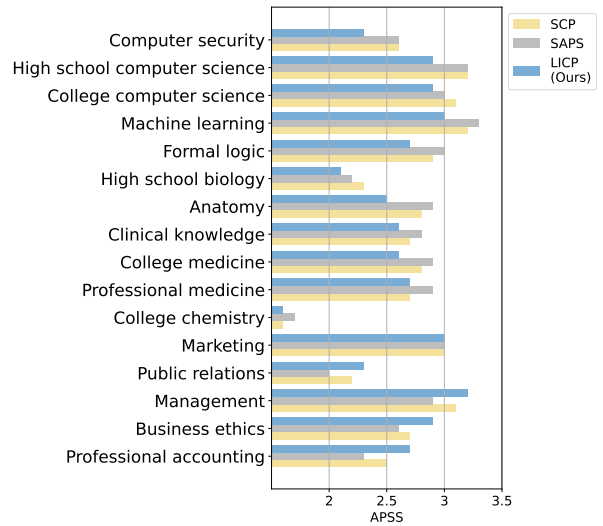


Figure 4: Results on MMLU for MCQ task, with the error rate of 0.2. Our method and baselines are applied individually to each of the 16 subjects.

that Vicuna-v1.5(7b) can only produce prediction sets with higher error rates compared to Llama-2 backbones. This is because Vicuna-v1.5(7b) is less powerful for these two datasets. This demonstrates that CP performance for LLMs is largely dependent on the performance of the backbone models.

Sampling Quantity The sampling quantity regulates the number and types of sampled responses acquired from LLMs, thereby influencing frequency, NE and SS. We vary the sampling quantity from 10 to 40 on the TriviaQA dataset using Llama-2-13b, incrementing by 5 each time. Results shown in Figure 6 suggest that a larger sampling quantity tends to present better performance w.r.t. efficiency. This is because, with a higher sampling quantity, the frequency notion more accurately represents response rankings. Of particular interest is that, at an error rate of 0.2, the sampling quantity of 15 exhibits inferior performance compared to the quantity of 10. We hypothesize it is because a sampling quantity

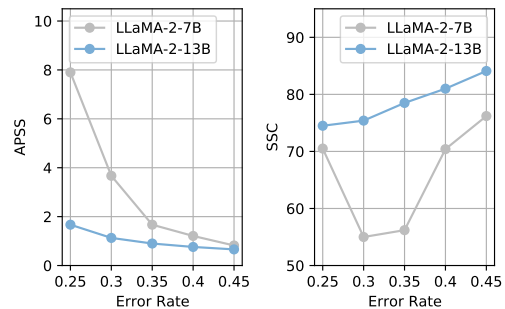


Figure 5: Results of the sensitivity analysis for different backbone models: Llama-2-7b and Llama-2-13b.

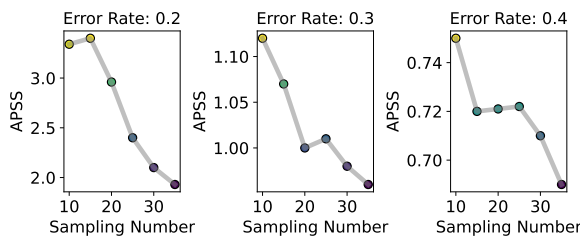


Figure 6: Sensitivity analysis of sampling quantity.

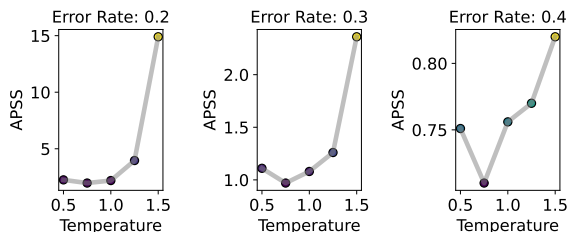


Figure 7: Sensitivity analysis of temperature.

of 15 remains insufficient to adequately represent rankings meanwhile introducing more non-robust randomness in responses. In addition, we observe a larger impact of the sampling quantity on APSS when a small error rate guarantee is required.

Temperature Scaling. The temperature (Hinton et al., 2015) in LLMs adjusts the randomness in generated outputs by scaling logits during the softmax operation. Higher temperatures boost the diversity of the output, which may further affect the performance of LofreeCP. In this experiment, we vary temperatures² (0.5, 0.75, 1.0, 1.25, and 1.5) in the Llama-2-13b model. Results for the TriviaQA dataset are presented in Figure 7. The smallest (best) APSS is observed at a temperature of 0.75. We observe an overall growing trend as the temperature increases from 0.75 to 1.50. This indicates that excessive diversity can result in uncertain and suboptimal predictions. The decline from 0.50 to 0.75 implies that too much determinism may hurt CP efficiency due to a lack of randomness and diversity. We also note a significant temperature influence on APSS when aiming for low error rates.

5 Related Work

Conformal Prediction for NLP. CP has already found diverse applications in NLP, e.g., text infilling and part-of-speech prediction Dey et al. (2021), sentiment analysis Maltoudoglou et al. (2020), and Automatic Speech Recognition Ernez et al. (2023). In the application of CP to LLMs, existing methods are predominantly logit-based. For instance,

²Temperature ranges between 0 and 2.

Kumar et al. (2023) apply standard CP (Vovk et al., 2005) to Llama-2-13b (Touvron et al., 2023) for the MCQ task by computing softmax scores of token logits for options to measure nonconformity. Similarly, Quach et al. (2023) extend CP to LLMs using the general risk control framework (Angelopoulos et al., 2021). However, recent studies have pointed out that relying solely on logits may be flawed due to hallucinations in LLMs (LeCun, 2023). Consequently, there is ongoing research aiming to reduce reliance on logits. Huang et al. (2023) propose to use the highest probability and replace other probabilities with weighted values. All these methods involve the utilization of logits.

Uncertainty Estimation in LLMs. Recent developments in LLMs have highlighted the importance of estimating their uncertainty. While there has been significant research on uncertainty in NLP (Van Landeghem et al., 2022; Ulmer et al., 2022), several methods exist to estimate the confidence of LLMs, including Deep Ensemble methods (Lakshminarayanan et al., 2017), Monte Carlo dropout (Gal and Ghahramani, 2016), Density-based estimation (Yoo et al., 2022), Confidence learning (DeVries and Taylor, 2018), as well as approaches based on logits. However, recent studies highlight concerns that LLMs may generate unfaithful and nonfactual content (Maynez et al., 2020). Additionally, logits of LLMs often exhibit overconfidence when producing incorrect answers, indicating that logits alone may not be entirely reliable for studying uncertainty (Desai and Durrett, 2020; Miao et al., 2021; Vasconcelos et al., 2023).

6 Conclusion

We study the critical problem of CP for API-only LLMs without logit-access. We propose a novel solution to define the nonconformity score function by leveraging uncertainty information from diverse sources. In particular, under a limited sampling budget, we first use the response frequency as the coarse-grained proxy of uncertainty levels. We then propose two fine-grained uncertainty notions (NE and SS) to further distinguish uncertainty at a nuanced level. Our proposed approach does not rely on model logits and can alleviate the known miscalibration issue when using logits. Experiments demonstrate the superior performance of our approach compared to logit-based and logit-free baselines. Our work opens up a new avenue to uncertainty estimation in LLMs without logit-access.

588 Limitations

589 Our approach encounters a common limitation of
590 open-ended Natural Language Generation (NLG)
591 tasks: the unbounded output space. In our work, we
592 address this challenge by sampling a fixed number
593 of times for every prompt from LLMs to achieve
594 a comprehensive output space, but we recognize
595 the potential for more effective and convincing ap-
596 proaches to handle this issue within the framework
597 of CP. Secondly, another future direction is to ex-
598 pand our CP method to non-exchangeability sce-
599 narios, particularly in NLG domains, where cali-
600 bration and test sets may not adhere strictly to
601 the assumption of being independent and identically
602 distributed (i.i.d.). Finally, due to financial
603 constraints, we do not evaluate our approach on
604 several proprietary LLMs (e.g., GPT 4) that allow
605 users to obtain token log probabilities. Thus future
606 work can validate our method on these models.

607 References

608 Anastasios Angelopoulos, Stephen Bates, Jitendra Ma-
609 lik, and Michael I Jordan. 2020. Uncertainty sets for
610 image classifiers using conformal prediction. *arXiv*
611 *preprint arXiv:2009.14193*.

612 Anastasios N Angelopoulos and Stephen Bates. 2021.
613 A gentle introduction to conformal prediction and
614 distribution-free uncertainty quantification. *arXiv*
615 *preprint arXiv:2107.07511*.

616 Anastasios N Angelopoulos, Stephen Bates, Em-
617 manuel J Candès, Michael I Jordan, and Lihua Lei.
618 2021. Learn then test: Calibrating predictive al-
619 gorithms to achieve risk control. *arXiv preprint*
620 *arXiv:2110.01052*.

621 Jonathan Berant, Andrew Chou, Roy Frostig, and Percy
622 Liang. 2013. Semantic parsing on freebase from
623 question-answer pairs. In *Proceedings of the 2013*
624 *conference on empirical methods in natural language*
625 *processing*, pages 1533–1544.

626 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
627 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
628 Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al.
629 2023. Vicuna: An open-source chatbot impressing
630 gpt-4 with 90%* chatgpt quality. See [https://vicuna.](https://vicuna.lmsys.org)
631 [lmsys.org](https://vicuna.lmsys.org) (accessed 14 April 2023).

632 Shrey Desai and Greg Durrett. 2020. [Calibration of](#)
633 [pre-trained transformers](#).

634 Terrance DeVries and Graham W. Taylor. 2018. [Learn-](#)
635 [ing confidence for out-of-distribution detection in](#)
636 [neural networks](#).

Neil Dey, Jing Ding, Jack Ferrell, Carolina Kapper, 637
Maxwell Lovig, Emiliano Planchon, and Jonathan P 638
Williams. 2021. [Conformal prediction for text infill-](#) 639
[ing and part-of-speech prediction](#). 640

Fares Ernez, Alexandre Arnold, Audrey Galametz, 641
Catherine Kobus, and Nawal Ould-Amer. 2023. [Ap-](#) 642
[plying the conformal prediction paradigm for the](#) 643
[uncertainty quantification of an end-to-end automatic](#) 644
[speech recognition model \(wav2vec 2.0\)](#). In *Proceed-* 645
ings of the Twelfth Symposium on Conformal and 646
Probabilistic Prediction with Applications, volume 647
204 of *Proceedings of Machine Learning Research*, 648
pages 16–35. PMLR. 649

Yarin Gal and Zoubin Ghahramani. 2016. [Dropout as](#) 650
[a bayesian approximation: Representing model un-](#) 651
[certainty in deep learning](#). In *Proceedings of The* 652
33rd International Conference on Machine Learn- 653
ing, volume 48 of *Proceedings of Machine Learning* 654
Research, pages 1050–1059, New York, New York, 655
USA. PMLR. 656

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, 657
Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 658
2020. Measuring massive multitask language under- 659
standing. *arXiv preprint arXiv:2009.03300*. 660

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. 661
Distilling the knowledge in a neural network. *arXiv* 662
preprint arXiv:1503.02531. 663

Jianguo Huang, Huajun Xi, Linjun Zhang, Huaxiu Yao, 664
Yue Qiu, and Hongxin Wei. 2023. Conformal pre- 665
diction for deep classifier via label ranking. *arXiv* 666
preprint arXiv:2310.06430. 667

Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke 668
Zettlemoyer. 2017. Triviaqa: A large scale distantly 669
supervised challenge dataset for reading comprehen- 670
sion. *arXiv preprint arXiv:1705.03551*. 671

Yuko Kato, David MJ Tax, and Marco Loog. 2023. A 672
review of nonconformity measures for conformal pre- 673
diction in regression. *Conformal and Probabilistic* 674
Prediction with Applications, pages 369–383. 675

Bhawesh Kumar, Charlie Lu, Gauri Gupta, Anil Palepu, 676
David Bellamy, Ramesh Raskar, and Andrew Beam. 677
2023. Conformal prediction with large language 678
models for multi-choice question answering. *arXiv* 679
preprint arXiv:2305.18404. 680

Balaji Lakshminarayanan, Alexander Pritzel, and 681
Charles Blundell. 2017. [Simple and scalable pre-](#) 682
[dictive uncertainty estimation using deep ensembles](#). 683
In *Advances in Neural Information Processing Sys-* 684
tems, volume 30. Curran Associates, Inc. 685

Y LeCun. 2023. Do large language models need sen- 686
sory grounding for meaning and understanding? In 687
Workshop on Philosophy of Deep Learning, NYU 688
Center for Mind, Brain, and Consciousness and the 689
Columbia Center for Science and Society. 690

691	Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	742
692	Jian-Guang Lou, and Weizhu Chen. 2022. On the	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	743
693	advance of making language models better reasoners.	Baptiste Rozière, Naman Goyal, Eric Hambro,	744
694	<i>arXiv preprint arXiv:2206.02336</i> .	Faisal Azhar, et al. 2023. Llama: Open and effi-	745
		cient foundation language models. <i>arXiv preprint</i>	746
695	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	<i>arXiv:2302.13971</i> .	747
696	Teaching models to express their uncertainty in		
697	words. <i>arXiv preprint arXiv:2205.14334</i> .	Dennis Ulmer, Jes Frellsen, and Christian Hardmeier.	748
		2022. Exploring predictive uncertainty and calibra-	749
698	Lysimachos Maltoudoglou, Andreas Paisios, and Harris	tion in NLP: A study on the impact of method & data	750
699	Papadopoulos. 2020. Bert-based conformal predictor	scarcity. In <i>Findings of the Association for Computa-</i>	751
700	for sentiment analysis. In <i>Proceedings of the Ninth</i>	<i>tional Linguistics: EMNLP 2022</i> , pages 2707–2735,	752
701	<i>Symposium on Conformal and Probabilistic Predic-</i>	Abu Dhabi, United Arab Emirates. Association for	753
702	<i>tion and Applications</i> , volume 128 of <i>Proceedings of</i>	Computational Linguistics.	754
703	<i>Machine Learning Research</i> , pages 269–284. PMLR.		
		Jordy Van Landeghem, Matthew Blaschko, Bertrand	755
704	James Manyika and Sissie Hsiao. 2023. An overview	Anckaert, and Marie-Francine Moens. 2022. Bench-	756
705	of bard: an early experiment with generative ai. <i>AI</i>	marking scalable predictive uncertainty in text classi-	757
706	<i>Google Static Documents</i> , 2.	fication. <i>IEEE Access</i> , 10:43703–43737.	758
707	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	Helena Vasconcelos, Gagan Bansal, Adam Fourney,	759
708	Ryan McDonald. 2020. On faithfulness and factu-	Q. Vera Liao, and Jennifer Wortman Vaughan. 2023.	760
709	ality in abstractive summarization. In <i>Proceedings</i>	Generation probabilities are not enough: Exploring	761
710	<i>of the 58th Annual Meeting of the Association for</i>	the effectiveness of uncertainty highlighting in ai-	762
711	<i>Computational Linguistics</i> , pages 1906–1919, On-	powered code completions.	763
712	line. Association for Computational Linguistics.		
		Vladimir Vovk, Alexander Gammerman, and Glenn	764
713	Mengqi Miao, Fandong Meng, Yijin Liu, Xiao-Hua	Shafer. 2005. <i>Algorithmic learning in a random</i>	765
714	Zhou, and Jie Zhou. 2021. Prevent the language	<i>world</i> , volume 29. Springer.	766
715	model from being overconfident in neural machine		
716	translation.	Fangxin Wang, Lu Cheng, Ruocheng Guo, Kay Liu, and	767
		Philip S Yu. 2023. Equal opportunity of coverage	768
717	Khanh Nguyen and Brendan O’Connor. 2015. Pos-	in fair regression. <i>Advances in Neural Information</i>	769
718	terior calibration and exploratory analysis for nat-	<i>Processing Systems</i> , 36.	770
719	ural language processing models. <i>arXiv preprint</i>		
720	<i>arXiv:1508.05154</i> .	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	771
		Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	772
721	OpenAI. 2023. GPT-4v(ision): technical work.	Denny Zhou. 2022. Self-consistency improves chain	773
		of thought reasoning in language models. <i>arXiv</i>	774
722	René Peinl and Johannes Wirth. 2023. Evaluation of	<i>preprint arXiv:2203.11171</i> .	775
723	medium-large language models at zero-shot closed		
724	book generative question answering. <i>arXiv preprint</i>	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	776
725	<i>arXiv:2305.11991</i> .	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	777
		et al. 2022. Chain-of-thought prompting elicits rea-	778
726	Victor Quach, Adam Fisch, Tal Schuster, Adam Yala,	soning in large language models. <i>Advances in Neural</i>	779
727	Jae Ho Sohn, Tommi S Jaakkola, and Regina Barzilay.	<i>Information Processing Systems</i> , 35:24824–24837.	780
728	2023. Conformal language modeling. <i>arXiv preprint</i>		
729	<i>arXiv:2306.10193</i> .	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert	781
		Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,	782
730	Mauricio Sadinle, Jing Lei, and Larry Wasserman. 2019.	Da Huang, Denny Zhou, et al. 2023. Larger language	783
731	Least ambiguous set-valued classifiers with bounded	models do in-context learning differently. <i>arXiv</i>	784
732	error levels. <i>Journal of the American Statistical As-</i>	<i>preprint arXiv:2303.03846</i> .	785
733	<i>sociation</i> , 114(525):223–234.		
		Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le,	786
734	Glenn Shafer and Vladimir Vovk. 2008. A tutorial on	Mohammad Norouzi, Wolfgang Macherey, Maxim	787
735	conformal prediction. <i>Journal of Machine Learning</i>	Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al.	788
736	<i>Research</i> , 9(3).	2016. Google’s neural machine translation system:	789
		Bridging the gap between human and machine trans-	790
737	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam	lation. <i>arXiv preprint arXiv:1609.08144</i> .	791
738	Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng,		
739	Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al.	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,	792
740	2022. Lamda: Language models for dialog applica-	Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin	793
741	tions. <i>arXiv preprint arXiv:2201.08239</i> .	Jiang. 2023. Wizardlm: Empowering large lan-	794
		guage models to follow complex instructions. <i>arXiv</i>	795
		<i>preprint arXiv:2304.12244</i> .	796

KiYoon Yoo, Jangho Kim, Jiho Jang, and Nojun Kwak. 2022. Detection of word adversarial examples in text classification: Benchmark and baseline via robust density estimation.

A Theoretical Proofs

A.1 Proof of Lemma 3.1

Proof. When N_{total} is sufficiently large, the Lindeberg-Lévy central limit theorem yields the following equation:

$$\frac{\frac{freq(Y_i)}{N_{total}} - p_i}{\sqrt{p_i(1-p_i)/N_{total}}} \sim N(0, 1),$$

From this, we conclude that

$$P \left\{ \left| \frac{\frac{freq(Y_i)}{N_{total}} - p_i}{\sqrt{p_i(1-p_i)/N_{total}}} \right| \leq u_{1-(1-\delta)/2} \right\} \geq \delta.$$

Approximately replacing p_i in the denominator with $\frac{freq(Y_i)}{N_{total}}$, we obtain

$$P \left\{ \left| \frac{\frac{freq(Y_i)}{N_{total}} - p_i}{\sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})/N_{total}}} \right| \leq u_{1-(1-\delta)/2} \right\} \geq \delta.$$

This equation is equivalent to

$$P \left\{ -u_{1-(1-\delta)/2} \leq \frac{\frac{freq(Y_i)}{N_{total}} - p_i}{\sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})/N_{total}}} \leq u_{1-(1-\delta)/2} \right\} \geq \delta.$$

We can then reformulate the above equation as:

$$\begin{aligned} & P \left\{ \frac{freq(Y_i)}{N_{total}} \right. \\ & \quad \left. - u_{1-(1-\delta)/2} \cdot \sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})} \right. \\ & \quad \leq p_i \leq \frac{freq(Y_i)}{N_{total}} \\ & \quad \left. + u_{1-(1-\delta)/2} \cdot \sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})} \right\} \\ & \geq \delta. \end{aligned}$$

In the left part of this equation, $\frac{freq(Y_i)}{N_{total}}$ represents the absolute frequency, p_i represents the desired estimated probability, and $u_{1-(1-\delta)/2}$

$\sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})}$ is the error term between them. Recall that we aim to ensure:

$$P \left\{ \left| \frac{freq(Y_i)}{N_{total}} - p_i \right| \leq \epsilon \right\} \geq \delta.$$

Therefore, we need to guarantee:

$$u_{1-(1-\delta)/2} \cdot \sqrt{\frac{freq(Y_i)}{N_{total}}(1 - \frac{freq(Y_i)}{N_{total}})} \cdot 2 \leq 2\epsilon.$$

This implies that we must control the error term to not exceed our predetermined estimation error. Note that the left part of this equation reaches its maximum value when $freq(Y_i) = \frac{1}{2}$. Hence, to achieve this, we only require:

$$\sqrt{\frac{1/4}{N_{total}}} \cdot u_{1-(1-\delta)/2} \cdot 2 \leq 2\epsilon.$$

This simplifies to

$$N_{total} \geq \left(\frac{u_{1-(1-\delta)/2}}{2\epsilon} \right)^2$$

□

A.2 Proof of Proposition 3.2

Proof. Let N denote the nonconformity measures of the calibration set $(X_i, Y_i)_{i=1, \dots, n}$, and let α_1 and α_2 be the desired error rates, where $\alpha_1 > \alpha_2$. As indicated in Step 2, we have $\hat{q}_1 \leq \hat{q}_2$. Given $C(X_{test}) = \{Y : N(X_{test}, Y) \leq \hat{q}\}$, it follows that $C_{1-\alpha_1}(X) \subseteq C_{1-\alpha_2}(X)$. Consequently, the nesting property, as defined in Equation 1, is satisfied. Therefore, Proposition 3.2 holds. □

B Implementation Details

B.1 Dataset

The TriviaQA benchmark (available at <https://nlp.cs.washington.edu/triviaqa/>) or can be accessed from Hugging Face at https://huggingface.co/datasets/trivia_qa) and the WebQuestions benchmark (available at worksheets.codalab.org) or can be accessed from Hugging Face at https://huggingface.co/datasets/web_questions) are employed for QA. Both datasets operate within a closed-book setting, where LLMs refrain from using supporting text when answering questions.

The MMLU benchmark (can be accessed from Hugging Face at <https://huggingface.co/>)

co/datasets/lukaemon/mmlu) is designed for MCQ, which covers 57 subjects across STEM, the humanities, the social sciences, and more. For our MCQ experiments, we leverage the dataset containing 16 subjects from the MMLU: computer security, high school computer science, college computer science, machine learning, formal logic, high school biology, anatomy, clinical knowledge, college medicine, professional medicine, college chemistry, marketing, public relations, management, business ethics, professional accounting.

For the TriviaQA dataset, we randomly select 10,000 question-answer pairs. Similarly, for the WebQuestions dataset, we randomly select 5,000 question-answer pairs. Regarding the MMLU dataset, we use all available data for each of the 16 subjects. Across all three datasets, we apply the same splitting strategy: 50% of the data serves as the calibration set, 25% as the validation set, and 25% as the test set for each trial.

B.2 Backbone LLMs

We utilize the Hugging Face API to access open-source LLMs in our experiments, including Llama-2-7B (accessible at huggingface.co/meta-llama/Llama-2-7b-hf), Llama-2-13B (accessible at huggingface.co/meta-llama/Llama-2-13b-hf), WizardLM-v1.2(13b) (accessible at huggingface.co/WizardLM/WizardLM-13B-V1.2), and Vicuna-v1.5(7b) (accessible at huggingface.co/lmsys/vicuna-7b-v1.5). Access to Llama-2-7b and Llama-2-13b requires requesting approval via the Meta website (<https://llama.meta.com/>). Upon approval, access to these resources will be granted.

B.3 Length-Normalization

We use length normalization (Wu et al., 2016) on logits to obtain response probability/likelihood:

$$p(x, y_k) = \exp\left(\frac{\log p_\theta(y_k|x)}{lp(y_k)}\right)$$

where

$$lp(y) = \frac{(5 + |y|)^{0.6}}{(5 + 1)^{0.6}}$$

B.4 Evaluation

We extract an answer by analyzing the text until we encounter the first line break, comma, or period. This implies that in the dataset, we will disre-

gard data whose answers contain line breaks, commas, or periods. Following this, we standardize the answers by converting them to lowercase, removing articles, punctuation, and duplicate whitespace. The generated answers are then evaluated using the exact match metric, where an answer is considered correct only if it exactly matches the provided answer. These guidelines align with those described in Quach et al. (2023).

For SSC, We focus exclusively on bins with a set size greater than 0 and a sample number exceeding 10% of the total test samples. This is because bins with a size of 0 and fewer samples lack reliability for coverage measurement.

B.5 LLMs Parameters

We employ the default Transformer generative LLMs parameters for our experiments, using default standard sampling with do_sample set to True, top_k set to 0, top_p set to 1, and Temperature set to 1, except when conducting model hyperparameter-tuning experiments. In such hyperparameter-tuning cases, we explicitly mention the parameters in main body of the paper.

B.6 Semantic Similarity

The measure of semantic similarity was established leveraging the FastText model available within the gensim package. The configuration parameters were carefully selected, defining a vector size of 200 and imposing a minimum count threshold of 1 to ensure robustness and inclusivity in the model’s representations.

B.7 Experiment trails

We conduct 50 trials for all experiments, then average the results to eliminate randomness during the calibration.

B.8 Error Rate Settings

We do not apply the same error rate settings across different models or datasets. This is because each model varies in its coverage ability for the same dataset. Likewise, the same model doesn’t possess identical coverage abilities for different datasets. Therefore, we adjust error rate settings for different combinations of model and dataset accordingly.

B.9 GPUs

We utilize six NVIDIA RTX 3090 graphics cards to support experiments.

1022	Q: Which prince is Queen Elizabeth II's	the beginning of WW2?	1072
1023	youngest son?	A: Germany	1073
1024	A: Edward	Q: Which countries border the US?	1074
1025	Q: When did the founder of Jehovah's	A: Canada	1075
1026	Witnesses say the world would end?	Q: Where is Rome, Italy located on a	1076
1027	A: 1914	map?	1077
1028	Q: Who found the remains of the Titanic?	A: Rome	1078
1029	A: Robert Ballard	Q: What is Nina Dobrev's nationality?	1079
1030	Q: Who was the only Spice Girl not to	A: Bulgaria	1080
1031	have a middle name?	Q: What country does Iceland belong to?	1081
1032	A: Posh Spice	A: Iceland	1082
1033	Q: What are the international	Q: What does Thai mean?	1083
1034	registration letters of a vehicle from	A: Language	1084
1035	Algeria?	Q: Who was Ishmael's mom?	1085
1036	A: DZ	A: Hagar	1086
1037	Q: How did Jock die in Dallas?	Q: What are the major cities in France?	1087
1038	A: Helicopter accident	A: Paris	1088
1039	Q: What star sign is Michael Caine?	Q: What city did Esther live in?	1089
1040	A: Pisces	A: Susa	1090
1041	Q: Who wrote the novel Evening Class?	Q: What sport do the Toronto Maple Leafs	1091
1042	A: Maeve Binchy	play?	1092
1043	Q: Which country does the airline Air	A: Ice Hockey	1093
1044	Pacific come from?	Q: What is Martin Cooper doing now?	1094
1045	A: Fiji	A: Inventor	1095
1046	Q: In which branch of the arts does	Q: What county is the city of Hampton,	1096
1047	Allegra Kent work?	VA in?	1097
1048	A: Ballet	A: Hampton	1098
1049	Q: Banting and Best pioneered the use of	Q: What county is Heathrow Airport in?	1099
1050	what?	A: London	1100
1051	A: Insulin	Q: What type of car does Michael Weston	1101
1052	Q: Who directed the movie La Dolce Vita?	drive?	1102
1053	A: Federico Fellini	A: Wishcraft	1103
1054	Q: Which country does the airline LACSA	Q: What was Tupac's name in Juice?	1104
1055	come from?	A: Bishop	1105
1056	A: Costa Rica	Q: Who does Maggie Grace play in Taken?	1106
1057	Q: Who directed 2001: A Space Odyssey?	A: Kim	1107
1058	A: Stanley Kubrick	Q: What style of music did Louis	1108
1059	Q: Which is the largest of the Japanese	Armstrong play?	1109
1060	Volcano Islands?	A: Jazz	1110
1061	A: Iwo Jima	Q: Where does Jackie French live?	1111
1062	Q: (Question)	A: Australia	1112
1063	A:	Q: Where is Jack Daniels factory?	1113
		A: Tennessee	1114
1064	C.2 Prompts of Webquestions	Q: What is Charles Darwin famous for?	1115
1065	We also use 32-shot question-answer pair prompts	A: Evolution	1116
1066	from the Webquestions train set.	Q: Where to visit in N. Ireland?	1117
1067	Answer these questions.	A: Antrim	1118
1068	Q: What country is the Grand Bahama	Q: What are dollars called in Spain?	1119
1069	Island in?	A: Peseta	1120
1070	A: Bahamas	Q: Who plays Meg in Family Guy?	1121
1071	Q: What two countries invaded Poland in	A: Mila Kunis	1122

1123 Q: Where did Martin Luther King get
 1124 shot?
 1125 A: Memphis
 1126 Q: What was Nelson Mandela’s religion?
 1127 A: Methodism
 1128 Q: Who will win the 2011 NHL Stanley
 1129 Cup?
 1130 A: Canada
 1131 Q: What is Henry Clay known for?
 1132 A: Lawyer
 1133 Q: What is the money of Spain called?
 1134 A: Euro
 1135 Q: Where are Sunbeam microwaves made?
 1136 A: Florida
 1137 Q: Where was Kennedy when he got shot?
 1138 A: Dallas
 1139 Q: Where did the Casey Anthony case take
 1140 place?
 1141 A: Orlando
 1142 Q: (Question)
 1143 A:

1144 C.3 Prompts of MMLU

1145 Each subject in MMLU uses similar prompts. We
 1146 take the high school biology as examples.

1147 Please engage in the multiple-choice
 1148 question-answering task. You should
 1149 generate the option (A, B, C, or D) you
 1150 think is right. Examples are provided.
 1151 (Select 8-shot randomly from other
 1152 subjects)

1153 This is a question from high school
 1154 biology.
 1155 A piece of potato is dropped into a
 1156 beaker of pure water. Which of the
 1157 following describes the activity after
 1158 the potato is immersed into the water?
 1159 (A) Water moves from the potato into the
 1160 surrounding water.
 1161 (B) Water moves from the surrounding
 1162 water into the potato.
 1163 (C) Potato cells plasmolyze.
 1164 (D) Solutes in the water move into the
 1165 potato.
 1166 The correct answer is option: B.

1167 You are the world’s best expert in high
 1168 school biology. Reason step-by-step and
 1169 answer the following question.
 1170 From the solubility rules, which of the
 1171 following is true?

(A) All chlorides, bromides, and iodides 1172
 are soluble 1173
 (B) All sulfates are soluble 1174
 (C) All hydroxides are soluble 1175
 (D) All ammonium-containing compounds 1176
 are soluble 1177
 The correct answer is option: 1178

1179 D Additional Results

1180 D.1 Ablation Study

Table 3: SCC Results of Ablation Study

Error Rate	0.20	0.25	0.30	0.35	0.40	0.45
Freq-Only	77.1	72.9	75.3	77.2	79.4	81.7
Freq + NE	78.8	74.0	76.8	77.9	80.2	83.3
Freq + SS	78.2	74.7	76.6	78.7	80.0	82.9
All (Ours)	79.2	74.3	76.5	78.6	81.5	84.0

Table 4: Portion of Concentration

Method	Portion of Concentration (%)
Freq-Only	66.6
Freq + NE	45.5
Freq + SS	59.8
All (Ours)	35.1

1181 D.2 Sensitivity Experiments

1182 More results regarding sampling quantity and tem-
 1183 perature sensitivity are included in Figures 8-9 due
 1184 to the page limit in the main body. 1184

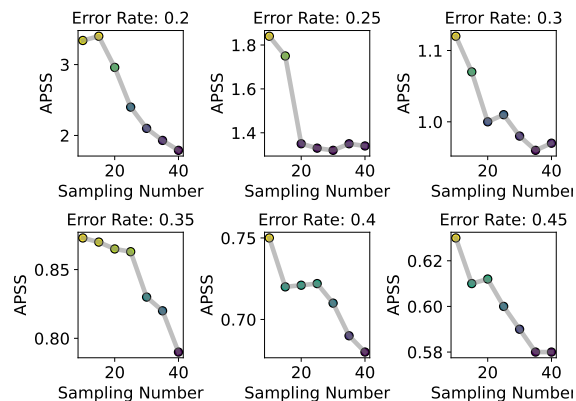


Figure 8: All results of the sensitivity analysis to variations in sampling quantity.

1185 D.3 Results for WizardLM-v1.2 (13B) and 1186 Vicuna-v1.5 (7B)

1187 To save on computation costs, we use float16 preci-
 1188 sion (half-precision) for experiments in this section. 1188
 1189 We use standard sampling with sampling quantity 1189

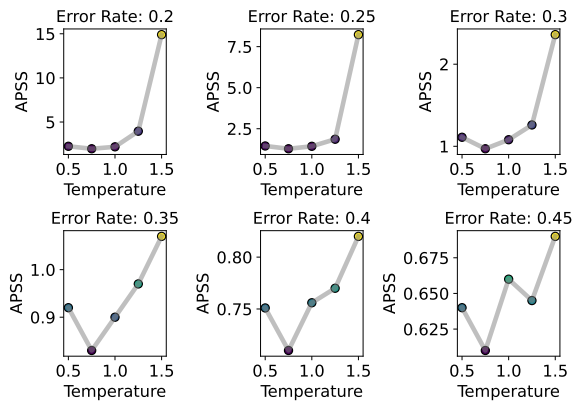


Figure 9: All results of the sensitivity analysis to variations in temperature.

1190 of 30. Results for TriviaQA are shown in Table 5,
 1191 for WebQuestions are shown in Table 7. Results
 1192 for TriviaQA are shown in Table 6, for WebQues-
 1193 tions are shown in Table 8. Results for WizardLM-
 1194 v1.2 (13B) and Vicuna-v1.5 (7B) consistently align
 1195 with the main body results, demonstrating that the
 1196 LofreeCP method mostly outperforms baselines.

Table 5: Results for TriviaQA using WizardLM-v1.2.

Methods	Logit-Access	Error Rate								
		0.25			0.3			0.35		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	75.1	68.7	3.19	71.0	65.8	2.56	66.4	63.3	1.84
CLM	✓	75.1	63.3	3.01	70.1	64.9	2.20	65.0	63.3	1.43
SCP	✓	75.4	57.9	3.29	70.1	62.2	2.15	65.2	56.4	1.68
SAPS	✓	75.1	70.6	3.83	70.1	53.2	2.30	65.1	54.9	1.37
First-K _{black}	✗	75.7	58.0	4.94	71.5	66.6	2.59	68.4	65.6	1.84
LofreeCP (Ours)	✗	75.1	68.0	4.07	70.0	67.7	1.92	65.1	70.1	1.27

Methods	Logit-Access	Error Rate								
		0.4			0.45			0.5		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	✗	✗	✗	55.2	56.0	0.99	✗	✗	✗
CLM	✓	60.1	65.3	1.25	55.1	69.1	0.92	50.1	71.3	0.81
SCP	✓	60.0	65.9	1.30	55.1	67.8	1.01	50.1	70.1	0.82
SAPS	✓	60.0	47.3	1.37	55.2	53.7	1.05	50.1	60.6	0.83
First-K _{black}	✗	✗	✗	✗	56.9	57.4	0.99	✗	✗	✗
LofreeCP (Ours)	✗	60.2	69.8	0.98	55.3	70.4	0.81	50.2	72.5	0.69

Table 6: Results for TriviaQA using Vicuna-v1.5.

Methods	Logit-Access	Error Rate								
		0.475			0.5			0.525		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	53.0	42.1	2.23	50.4	42.4	1.63	✗	✗	✗
CLM	✓	52.5	45.1	2.60	50.1	45.5	1.39	47.5	47.7	1.21
SCP	✓	52.6	39.0	2.66	50.0	40.5	1.43	47.9	49.3	1.14
SAPS	✓	52.7	40.1	2.30	50.3	48.8	1.59	47.5	45.6	1.24
First-K _{black}	✗	53.4	44.1	2.75	50.9	42.3	1.62	✗	✗	✗
LofreeCP (Ours)	✗	52.5	39.3	2.27	50.0	39.1	1.33	47.6	50.1	1.12

Methods	Logit-Access	Error Rate								
		0.4			0.45			0.5		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	45.0	46.7	0.99	✗	✗	✗	✗	✗	✗
CLM	✓	45.2	50.7	1.01	42.5	50.6	0.85	40.1	56.2	0.83
SCP	✓	45.4	52.4	0.96	42.6	48.6	0.85	40.5	52.0	0.76
SAPS	✓	45.0	46.2	1.04	42.6	50.8	0.84	40.1	57.9	0.75
First-K _{black}	✗	✗	✗	✗	44.6	46.2	0.97	✗	✗	✗
LofreeCP (Ours)	✗	45.1	55.3	0.96	42.7	58.0	0.82	40.2	58.5	0.73

Table 7: Results for WebQuestions using WizardLM-v1.2.

Methods	Logit-Access	Error rate											
		0.45			0.5			0.55			0.6		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	55.5	42.5	3.40	53.0	40.6	2.70	49.1	39.0	1.91	✗	✗	✗
CLM	✓	55.1	52.3	3.02	50.2	40.1	2.01	45.2	28.6	1.58	40.4	31.2	1.19
SCP	✓	55.2	45.9	3.63	50.1	40.8	2.04	45.0	37.1	1.55	40.2	47.8	1.04
SAPS	✓	55.0	45.7	3.38	50.1	41.1	2.15	45.2	28.6	1.58	40.4	31.2	1.19
First-K _{black}	✗	56.7	43.6	3.40	50.9	45.0	1.91	✗	✗	✗	41.4	41.1	1.00
LofreeCP (Ours)	✗	55.0	45.3	2.87	50.0	46.5	1.88	45.1	49.9	1.18	40.1	51.7	0.82

Table 8: Results for WebQuestions using Vicuna-v1.5.

Methods	Logit-Access	Error rate											
		0.575			0.6			0.625			0.65		
		ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow	ECR	SSC \uparrow	APSS \downarrow
First-K _{white}	✓	43.2	23.8	1.99	41.7	26.9	1.57	✗	✗	✗	36.6	36.6	1.00
CLM	✓	42.5	32.3	1.88	40.1	36.2	1.32	37.6	38.2	1.08	35.0	41.8	0.83
SCP	✓	42.6	31.1	1.91	40.1	34.4	1.28	38.2	37.3	1.06	35.2	43.7	0.87
SAPS	✓	42.5	32.3	1.88	40.1	36.2	1.32	37.6	38.2	1.08	35.0	41.8	0.83
First-K _{black}	✗	43.7	25.9	2.01	40.9	25.5	1.57	✗	✗	✗	36.8	36.8	1.00
LofreeCP (Ours)	✗	42.5	32.4	1.73	40.1	36.7	1.22	37.5	39.6	0.97	35.0	39.3	0.81