

RETHINKING UNCERTAINTY ESTIMATION IN LLMs: A PRINCIPLED SINGLE-SEQUENCE MEASURE

Lukas Aichberger^{1,*}, Kajetan Schweighofer^{1,*}, Sepp Hochreiter^{1,2}

¹ ELLIS Unit Linz and LIT AI Lab, Institute for Machine Learning,
Johannes Kepler University Linz, Austria

² NXAI GmbH, Linz, Austria

* Joint first authors

{aichberger, schweighofer, hochreit}@ml.jku.at

ABSTRACT

Large Language Models (LLMs) are increasingly employed in real-world applications, driving the need to evaluate the trustworthiness of their generated text. To this end, reliable uncertainty estimation is essential. Leading uncertainty estimation methods generate and analyze multiple output sequences, which is computationally expensive and impractical at scale. In this work, we inspect the theoretical foundations of these methods and explore new directions to enhance computational efficiency. Building on the framework of proper scoring rules, we find that the negative log-likelihood of the most likely output sequence constitutes a theoretically principled uncertainty measure. To approximate this alternative measure, we propose G-NLL, obtained using a single output sequence from greedy decoding. This approach streamlines uncertainty estimation while preserving theoretical rigor. Empirical results demonstrate that G-NLL achieves state-of-the-art performance across various scenarios. Our work lays the theoretical foundation for efficient and reliable uncertainty estimation in natural language generation, challenging the necessity of the prevalent methods that are more complex and resource-intensive.

1 INTRODUCTION

Despite advances in natural language generation (NLG), determining the trustworthiness of generated text remains challenging. Addressing this requires reliably estimating the uncertainty a language model has regarding its generated text. Although a low level of uncertainty does not guarantee factual correctness, particularly when the generated text is based on consistent but inaccurate training data, uncertainty estimates remain a reliable indicator of errors at present (Farquhar et al., 2024).

Assessing predictive uncertainty in Large Language Models (LLMs) is challenging due to their stochastic, autoregressive nature. Each token is selected probabilistically, leading to diverse outputs for the same input. Furthermore, the vast space of possible sequences is computationally intractable. Common uncertainty estimation methods thus rely on expectations over output distributions, such as sequence entropy (Malinin and Gales, 2021; Kuhn et al., 2023; Duan et al., 2024; Farquhar et al., 2024), which in turn requires sampling multiple output sequences. However, this is computationally costly due to the large number of model parameters. As a result, only a small subset of outputs is sampled in practice. However, differences between sampled sequences do not always indicate uncertainty, as they may vary lexically while remaining semantically similar. Some methods use inference models to assess semantics (Kuhn et al., 2023; Aichberger et al., 2025; Farquhar et al., 2024), improving uncertainty estimates but adding complexity and additional computation. These challenges make large-scale uncertainty estimation impractical for real-world applications.

Efficient uncertainty estimation methods are needed to ensure the trustworthiness of the language model’s answer without imposing excessive computational demands. To address this need, we theoretically motivate that uncertainty measures require only a single output sequence. We do so by building on insights from the framework of proper scoring rules (Gneiting and Raftery, 2007) that has recently been investigated for uncertainty estimation in the standard classification setting (Kotelevskii et al., 2025; Hofman et al., 2024). Specifically, we extend proper scoring rules for uncertainty

estimation to NLG and explore the zero-one score as an alternative to the prevalent logarithmic score. The resulting uncertainty measure is straightforward: it is the negative log-likelihood of the most likely output sequence. Importantly, the measure is based on the most likely output sequence, rather than an arbitrary one. However, determining the most likely output sequence remains computationally expensive. Therefore, we propose G-NLL, an approximation that is maximally efficient and minimizes algorithmic complexity, while maintaining state-of-the-art performance.

In parallel, recent work has considered the likelihood of a single output sequence for uncertainty estimation in NLG (Fadееva et al., 2023; 2024; Vazhentsev et al., 2024; Abbasi-Yadkori et al., 2024; Vashurin et al., 2025a;b). Notably, Fadееva et al. (2023) introduce the maximum sequence probability (MSP) as a baseline, which corresponds to the negative log-likelihood of the most likely output sequence. However, their ad hoc formulation does not provide a theoretical justification or a discussion on how to best approximate it. More generally, prior work has not properly characterized the MSP as a measure for uncertainty in NLG. As a result, prior work often relies on unfavorable length-normalization (Qiu and Miikkulainen, 2024; Chen et al., 2024; Bakman et al., 2024; Yaldiz et al., 2025), does not focus on the *most likely* output sequence, or even overlooks single-sequence measures entirely as baselines (Malinin and Gales, 2018; Manakul et al., 2023; Kuhn et al., 2023; Farquhar et al., 2024; Nikitin et al., 2024; Kossen et al., 2024; Duan et al., 2024).

To close this gap, we are the first to provide a theoretical justification for the MSP as a principled, single-sequence measure of uncertainty in NLG. We examine the properties of this alternative uncertainty measure and compare it to established measures from both theoretical and empirical perspectives. Additionally, we propose G-NLL as an efficient approximation of the MSP and analyze why sampling-based or length-normalized alternatives are detrimental to its approximation quality. Our experiments demonstrate that G-NLL matches and even exceeds the estimation quality of established methods across various model classes, model sizes, training stages, tasks, datasets, and evaluation metrics. By maintaining theoretical rigor, G-NLL offers an effective and scalable approach to uncertainty estimation in NLG. It serves not only as a strong baseline for future work, but also as a practical solution for deploying uncertainty estimation in LLMs across real-world applications.

To summarize, our main contributions are as follows:

- We derive the negative log-likelihood of the most likely output sequence (i.e., the MSP) as a measure of uncertainty in NLG, building on established principles from uncertainty estimation theory and proper scoring rules.
- We provide a theoretical analysis of the MSP and existing uncertainty measures, and show that estimating this single-sequence measure possesses desirable properties for practical scenarios.
- We propose G-NLL as an effective approximation of the MSP, and demonstrate that it empirically outperforms state-of-the-art methods while significantly reducing computational costs.

2 PREDICTIVE UNCERTAINTY IN NLG

We begin by reviewing language modeling to formalize predictive uncertainty in NLG. Sec. 2.1 introduces proper scoring rules and their connection to predictive uncertainty, and Sec. 2.2 revisits established measures within this framework. Sec. 2.3 then introduces the maximum sequence probability and its approximation G-NLL under an alternative scoring rule.

Preliminaries. We assume a fixed training dataset $\mathcal{D} = \{s_i\}_{i=1}^N$ consisting of N token sequences $s = (s_1, \dots, s_T)$ where individual tokens $s_t \in \mathcal{V}$ are from a given vocabulary \mathcal{V} . Each token at step t is assumed to be sampled according to the predictive distribution $p(s_t | s_{<t}, w^*)$, conditioned on the sequence of preceding tokens $s_{<t}$ and the true (but unknown) language model parameters w^* . We assume that the given model class can theoretically represent the true predictive distribution, a common and usually necessary assumption (Hüllermeier and Waegeman, 2021). The likelihood of some model parameters \tilde{w} matching w^* is given by the posterior $p(\tilde{w} | \mathcal{D}) = p(\mathcal{D} | \tilde{w})p(\tilde{w})/p(\mathcal{D})$.

The input to a given language model parameterized by w is a sequence $x = (x_1, \dots, x_M)$ and the output is a sequence $y = (y_1, \dots, y_T) \in \mathcal{Y}_T$, with $x, y \in \mathcal{V}$ and \mathcal{Y}_T being the set of all possible output sequences with sequence length T . The likelihood of a token $y_t \in y$ being generated by the language model is conditioned on both the input sequence and all previously generated tokens, denoted as $p(y_t | x, y_{<t}, w)$. The likelihood of output sequences $y \in \mathcal{Y}_T$ being generated by

the language model is then the product of the individual token probabilities, which is denoted as $p(\mathbf{y} | \mathbf{x}, \mathbf{w}) = \prod_{t=1}^T p(y_t | \mathbf{x}, \mathbf{y}_{<t}, \mathbf{w})$ (Sutskever et al., 2014), while the heuristic length-normalized variant is $\bar{p}(\mathbf{y} | \mathbf{x}, \mathbf{w}) = \exp \left\{ \frac{1}{T} \sum_{t=1}^T \log p(y_t | \mathbf{x}, \mathbf{y}_{<t}, \mathbf{w}) \right\}$ (Malinin and Gales, 2021).

Calculating the likelihood that a specific output sequence is generated by the language model parameterized by \mathbf{w} is straightforward. The language model directly provides the token likelihoods for a given input sequence. However, determining the full probability distribution on all possible output sequences is considerably more challenging, since the size of \mathcal{Y}_T increases exponentially with the sequence length. The computational complexity of evaluating all possible output sequences increases with $\mathcal{O}(|\mathcal{V}|^T)$. Since the vocabulary sizes $|\mathcal{V}|$ of modern LLMs are well over a hundred thousand tokens, this distribution becomes intractable to determine, even for relatively short maximal sequence lengths T (Dubey et al., 2024).

2.1 PROPER SCORING RULES AND THE RELATION TO UNCERTAINTY MEASURES IN NLG

We next give an introduction to proper scoring rules and discuss how they give rise to uncertainty measures. For more details, in the standard classification setting, we refer to Hofman et al. (2024); Kotelevskii et al. (2025). Proper scoring rules are a class of functions that evaluate the quality of probabilistic predictions by assigning a numerical score based on the predictive distribution and actual observations (Gneiting and Raftery, 2007). In particular, a proper scoring rule is an extended real-valued function $S : \mathcal{P} \times \mathcal{Y}_T \rightarrow [-\infty, \infty]$, such that $S(p, \mathbf{y})$ is \mathcal{P} -quasi-integrable over a convex class of probability measures \mathcal{P} . A scoring rule is called proper relative to \mathcal{P} if the expected score is minimized when the evaluated distribution $p \in \mathcal{P}$ coincides with the distribution from which the outcomes $\mathbf{y} \in \mathcal{Y}_T$ are sampled from, and it is called strictly proper if this minimum is unique. In the context of uncertainty estimation in NLG, proper scoring rules assign a numerical value that reflects how much probability the predictive distribution of the *true* model $p(\mathbf{y} | \mathbf{x}, \mathbf{w}^*)$ places on an observed output sequence \mathbf{y}' , denoted as $S(p(\mathbf{y} | \mathbf{x}, \mathbf{w}^*), \mathbf{y}')$. Here, the output sequence \mathbf{y}' is an independent notational copy of \mathbf{y} for clarity, where necessary. In the following, scoring rules are used in the loss convention, so larger expected scores correspond to higher predictive uncertainty.

To obtain concrete uncertainty measures, we need to make two specific assumptions (Schweighofer et al., 2025). First, we have to define the predictive distribution used to sample output sequences. Following Aichberger et al. (2025), we use a single, *given* language model with parameters \mathbf{w} to sample output sequences $\mathbf{y}' \sim p(\mathbf{y}' | \mathbf{x}, \mathbf{w})$. This assumption is also used in other works (Kuhn et al., 2023; Fadeeva et al., 2024; Farquhar et al., 2024) and is intuitively reasonable, since our main concern is the uncertainty of the output of a specific language model. Thus, we consider the expected score for possible output sequences under the predictive distribution of the given language model, denoted as $E_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})}[S(p(\mathbf{y} | \mathbf{x}, \mathbf{w}^*), \mathbf{y}')]$. This quantifies how well the predictive distribution of the given language model aligns with the true predictive distribution, capturing predictive uncertainty. Second, we have to define how the true model is approximated. We consider a Bayesian approximation of the true model, i.e., each possible language model $\tilde{\mathbf{w}}$ according to its posterior probability $p(\tilde{\mathbf{w}} | \mathcal{D})$ (Schweighofer et al., 2023b;a). Thus, we perform a posterior expectation over the expected score $E_{p(\tilde{\mathbf{w}} | \mathcal{D})}[E_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})}[S(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}')]]$, which can be additively decomposed into an entropy and a divergence term (Gneiting and Raftery, 2007; Kull and Flach, 2015):

$$\underbrace{E_{p(\tilde{\mathbf{w}} | \mathcal{D})} [E_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})} [S(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}')]]]}_{\text{expected score}} = \tag{1}$$

$$\underbrace{E_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})} [S(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}')]]}_{\text{entropy term}} + \underbrace{E_{p(\tilde{\mathbf{w}} | \mathcal{D})} [E_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})} [S(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}')] - S(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}')]]}_{\text{divergence term}}.$$

The expected score over possible output sequences \mathbf{y}' and language models $\tilde{\mathbf{w}}$ captures the *total* uncertainty of the *given* language model. The entropy term reflects *aleatoric* uncertainty, which quantifies the inherent stochasticity of generating output sequences with a given language model (Aichberger et al., 2025). The divergence term reflects *epistemic* uncertainty, which quantifies the uncertainty due to lack of knowledge about the true language model parameters arising from limited data (Houlsby et al., 2011; Gal, 2016; Malinin, 2019; Hüllermeier and Waegeman, 2021).

The remaining degree of freedom is the choice of a proper scoring rule $S(\cdot, \cdot)$, which determines the concrete form of the uncertainty measures. In the following, we consider two canonical proper scoring rules: the logarithmic score $S_{\log}(\cdot, \cdot)$ in Sec. 2.2 and the zero-one score $S_{0-1}(\cdot, \cdot)$ in Sec. 2.3.

2.2 ESTABLISHED UNCERTAINTY MEASURES IN NLG BASED ON THE LOGARITHMIC SCORE

The logarithmic score is typically assumed implicitly in both the standard classification (Houlsby et al., 2011; Gal, 2016) and the NLG setting (Malinin and Gales, 2021; Kuhn et al., 2023) to derive uncertainty measures. This is due to the grounding of the resulting uncertainty measures in information theory (Lahlou et al., 2023; Gruber and Buettner, 2023; Hofman et al., 2024; Kotelevskii et al., 2025). In the context of NLG, the logarithmic score considers the negative log-likelihood of a generated output sequence \mathbf{y}' :

$$S_{\log}(p(\mathbf{y} | \mathbf{x}, \cdot), \mathbf{y}') = -\log p(\mathbf{y} = \mathbf{y}' | \mathbf{x}, \cdot). \quad (2)$$

Substituting the logarithmic score into Eq. (1) results in the total uncertainty being the cross-entropy $\text{CE}(\cdot ; \cdot)$ between the output sequence distribution of the given model and those induced by models integrated over the entire posterior $p(\tilde{\mathbf{w}} | \mathcal{D})$ (Aichberger et al., 2025) (see Apx. A.1 for details):

$$\underbrace{E_{p(\tilde{\mathbf{w}}|\mathcal{D})}[\text{CE}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}); p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{total uncertainty}} = \underbrace{H(p(\mathbf{y} | \mathbf{x}, \mathbf{w}))}_{\text{aleatoric uncertainty}} + \underbrace{E_{p(\tilde{\mathbf{w}}|\mathcal{D})}[\text{KL}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) || p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{epistemic uncertainty}}. \quad (3)$$

The epistemic uncertainty is captured by a posterior expectation of the Kullback-Leibler divergence $\text{KL}(\cdot || \cdot)$ between the output sequence distribution of the given model and those induced by models across the entire posterior. This requires considering every possible parameterization of the model. Since modern LLMs have billions of parameters (Radford et al., 2018; Zhang et al., 2022; Dubey et al., 2024; Zuo et al., 2024), the epistemic uncertainty is challenging to estimate.

Thus, current work usually only focuses on the aleatoric uncertainty, captured by the Shannon entropy $H(\cdot)$ of the output sequence distribution of the given language model (Kuhn et al., 2023; Aichberger et al., 2025; Farquhar et al., 2024). Computing this output sequence distribution still requires considering the entire set of possible output sequences \mathcal{Y}_T , which is intractable and has to be approximated, as discussed in the following.

Predictive Entropy. The aleatoric uncertainty under a given language model is the entropy of the output sequence distribution $p(\mathbf{y} | \mathbf{x}, \mathbf{w})$, commonly referred to as predictive entropy (PE) (Malinin and Gales, 2021). Intuitively, high PE implies that the language model is likely to generate different output sequences from the same input sequence, indicating high uncertainty. PE is generally estimated by Monte Carlo (MC) sampling of output sequences (Malinin and Gales, 2021):

$$\begin{aligned} H(p(\mathbf{y} | \mathbf{x}, \mathbf{w})) &= E_{p(\mathbf{y}|\mathbf{x},\mathbf{w})}[-\log p(\mathbf{y} | \mathbf{x}, \mathbf{w})] \\ &\approx \frac{1}{N} \sum_{n=1}^N -\log p(\mathbf{y}^n | \mathbf{x}, \mathbf{w}), \quad \mathbf{y}^n \sim p(\mathbf{y} | \mathbf{x}, \mathbf{w}). \end{aligned} \quad (4)$$

Semantic Entropy. Semantic entropy (SE) (Kuhn et al., 2023; Farquhar et al., 2024) is based on the fact that output sequences may be different on a token level but equivalent on a semantic level. In such cases, the PE can be misleading, as it indicates high uncertainty even when different output sequences have the same semantic meaning. Thus, instead of the entropy of the output sequence distribution, the entropy of the semantic cluster distribution is considered, denoted as $p(c | \mathbf{x}, \mathbf{w}) = \sum_{\mathcal{Y}_T} p(c | \mathbf{x}, \mathbf{y}, \mathbf{w}) p(\mathbf{y} | \mathbf{x}, \mathbf{w})$. The probability of an output sequence belonging to a semantic cluster is usually approximated with a separate natural language inference model. SE thus measures uncertainty about the semantics of output sequences and is defined as

$$\begin{aligned} H(p(c | \mathbf{x}, \mathbf{w})) &= E_{p(c|\mathbf{x},\mathbf{w})}[-\log p(c | \mathbf{x}, \mathbf{w})] \\ &\approx \frac{1}{N} \sum_{n=1}^N -\log p(c^n | \mathbf{x}, \mathbf{w}), \quad c^n \sim p(c | \mathbf{x}, \mathbf{w}). \end{aligned} \quad (5)$$

For details on how to construct a tractable approximation of SE, we refer to Aichberger et al. (2025).

Discussion. In general, each of these uncertainty measures is based on the logarithmic score and considers the distribution over all possible output sequences $p(\mathbf{y} | \mathbf{x}, \mathbf{w})$, which is defined over the entire set of possible output sequences \mathcal{Y}_T . Approximating expectations over this distribution requires sampling multiple output sequences from \mathcal{Y}_T , which is computationally expensive. In the following, we show that this requirement can be eliminated when considering an alternative proper scoring rule.

2.3 NEW UNCERTAINTY MEASURES IN NLG BASED ON THE ZERO-ONE SCORE

In principle, any proper scoring rule may be used to derive viable uncertainty measures (Kotelevskii et al., 2025; Hofman et al., 2024). Thus, we explore measuring predictive uncertainty based on the zero-one score instead of the logarithmic score. Although it has been considered in the standard classification setting, to the best of our knowledge, the zero-one score has not yet been considered as a proper scoring rule for deriving uncertainty measures in NLG. The zero-one score only considers the predictive distribution for the *most likely output sequence*:

$$S_{0-1}(p(\mathbf{y} | \mathbf{x}, \cdot), \mathbf{y}') = 1 - \mathbb{1}\{\mathbf{y}' = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}_T} p(\mathbf{y} | \mathbf{x}, \cdot)\}. \quad (6)$$

Here, $\mathbb{1}\{\cdot\}$ denotes the indicator function, taking value 1 if its argument is true and 0 otherwise. Substituting the zero-one score into Eq. (1) results in the total uncertainty being the expected confidence that the given model assigns to the most likely output sequences of models across the entire posterior $p(\tilde{\mathbf{w}} | \mathcal{D})$ (see Apx. A.2 for details):

$$\underbrace{E_{p(\tilde{\mathbf{w}} | \mathcal{D})}[1 - p(\mathbf{y} = \tilde{\mathbf{y}}^* | \mathbf{x}, \mathbf{w})]}_{\text{total uncertainty}} = \underbrace{1 - p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w})}_{\text{aleatoric uncertainty}} + \underbrace{p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w}) - E_{p(\tilde{\mathbf{w}} | \mathcal{D})}[p(\mathbf{y} = \tilde{\mathbf{y}}^* | \mathbf{x}, \mathbf{w})]}_{\text{epistemic uncertainty}}, \quad (7)$$

where $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}_T} p(\mathbf{y} | \mathbf{x}, \mathbf{w})$ denotes the most likely output sequence under the given language model and $\tilde{\mathbf{y}}^* = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}_T} p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}})$ the most likely output sequence under any language model parametrized by $\tilde{\mathbf{w}}$.

As in Eq. (3), the epistemic uncertainty is a posterior expectation that remains challenging to estimate. Following prior work, we again focus on the aleatoric uncertainty, which considers the likelihood of the most likely output sequence under the given language model. It turns out that the aleatoric uncertainty is equivalent to the maximum sequence probability (MSP) previously introduced in an ad hoc manner (Fadeeva et al., 2023; 2024).

While the aleatoric uncertainty derived from the logarithmic score requires approximating an expectation over the entire output sequence distribution by sampling multiple output sequences (see Eq. (4) and Eq. (5)), the one derived from the zero-one score only requires approximating the most likely output sequence under the given language model (see Eq. (7)). This distinction is crucial, as sampling multiple output sequences is computationally expensive, while approximating the most likely output sequence aligns directly with standard inference techniques widely used in LLMs.

For numerical stability, it is preferable to consider the logarithm of the aleatoric uncertainty in Eq. (7) and omit the constant offset of one, as it does not affect the relative ordering of uncertainty estimates. This yields the negative log-likelihood (NLL) of the most likely output sequence, which we denote as

$$\begin{aligned} M(p(\mathbf{y} | \mathbf{x}, \mathbf{w})) &:= - \max_{\mathbf{y} \in \mathcal{Y}_T} \log p(\mathbf{y} | \mathbf{x}, \mathbf{w}) = - \max_{(y_1, \dots, y_T) \in \mathcal{Y}_T} \sum_t^T \log p(y_t | \mathbf{y}_{<t}, \mathbf{x}, \mathbf{w}) \quad (8) \\ &\propto \underbrace{1 - p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w})}_{\text{aleatoric uncertainty}} = 1 - \max_{(y_1, \dots, y_T) \in \mathcal{Y}_T} \prod_t^T p(y_t | \mathbf{y}_{<t}, \mathbf{x}, \mathbf{w}). \end{aligned}$$

We note that this single-sequence measure is also known as the min-entropy, the most conservative of the Rényi-entropies (Rényi, 1961), which is a lower bound on PE in Eq. (4).

With this, the task of uncertainty estimation reduces to finding the most likely output sequence. However, the search space is the set of all possible output sequences \mathcal{Y}_T , so the exact computation remains intractable. To make the computation tractable, we propose replacing the full sequence maximization with a tokenwise maximization by moving the max operator inside the summation. This yields a very efficient approximation, since it simply corresponds to standard greedy decoding with the given language model. Thus, we can approximate the NLL of the most likely output sequence in Eq. (8) using greedy decoding (G), which we formally define as

$$\text{G-NLL} := - \sum_{t=1}^T \log \left(\max_{y_t \in \mathcal{V}} p(y_t | \mathbf{x}, \mathbf{y}_{<t}, \mathbf{w}) \right) \approx M(p(\mathbf{y} | \mathbf{x}, \mathbf{w})). \quad (9)$$

Discussion. Although uncertainty measures based on the logarithmic score could, in principle, perform well if the distribution over output sequences $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ — or equivalently the distribution over semantic clusters $p(c \mid \mathbf{x}, \mathbf{w})$ — were tractable (as in standard classification tasks), these distributions are intractable for NLG due to the large vocabulary size and the auto-regressive nature of LLMs. As a result, sampling-based methods often yield crude approximations, suffering from computational costs and sampling variability. In contrast, G-NLL offers a principled alternative while eliminating the need for extensive sampling, making it highly practical and straightforward, while maintaining theoretical rigor through its foundation in proper scoring rules.

In general, uncertainty measures in NLG can be defined by a predictive distribution and a proper scoring rule: under the logarithmic score, the aleatoric uncertainty evaluated on $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ corresponds to PE (Eq. (4)), while evaluated on $p(c \mid \mathbf{x}, \mathbf{w})$ it corresponds to SE (Eq. (5)). Analogously, under the zero-one score, the aleatoric uncertainty evaluated on $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ corresponds to the maximum sequence probability (MSP), while evaluated on $p(c \mid \mathbf{x}, \mathbf{w})$ it yields the semantic counterpart of MSP, which we term the maximum cluster probability (MCP). G-NLL approximates MSP from a single output sequence, giving it a computational advantage over PE, whereas approximating MCP requires multiple output sequences, offering no such advantage over SE. We discuss this further in Apx. A.2. In the remainder of this paper, we focus on MSP and defer a detailed investigation of MCP to future work.

3 THEORETICAL ANALYSIS

Before empirically validating G-NLL in Sec. 5, we analyze the sample-complexity of approximating $H(\cdot)$ and $M(\cdot)$, showing that approximating the latter as done with G-NLL is desirable in the setting of LLMs. We consider a probability distribution of output sequences $p(\mathbf{y})$, dropping the dependency on \mathbf{x} and \mathbf{w} in this section for brevity. To incorporate importance sampling schemes, we consider access to a proposal distribution $q(\mathbf{y})$, e.g., by temperature-sampling from the LLM. We are interested in estimating $H(p(\mathbf{y}))$ and $M(p(\mathbf{y}))$ with their respective MC estimators:

$$\begin{aligned} M(p(\mathbf{y})) &= -\max_{\mathbf{y}} \log p(\mathbf{y}) & \hat{M}(p(\mathbf{y})) &= -\max \{ \log p(\mathbf{y}^n) \}_{n=1}^N \\ H(p(\mathbf{y})) &= -\mathbb{E}_{q(\mathbf{y})} \left[\frac{p(\mathbf{y})}{q(\mathbf{y})} \log p(\mathbf{y}) \right] & \hat{H}(p(\mathbf{y})) &= -\frac{1}{N} \sum_{n=1}^N \frac{p(\mathbf{y}^n)}{q(\mathbf{y}^n)} \log p(\mathbf{y}^n) \end{aligned}$$

where $\mathbf{y}^n \sim q(\mathbf{y})$. Furthermore, we assume boundedness, i.e., there exist constants a, b such that $a \leq \log p(\mathbf{y}) \leq b$ due to logit clipping and non-zero softmax temperatures. Then the following holds:

Theorem 1. *Let $\log p(\mathbf{y}) \in [a, b]$, $\forall \mathbf{y} \in \mathcal{Y}_T$ and $\epsilon > 0$. Then, with probability $1 - \delta$, we get:*

1. *Sample-Complexity Bound for Maximum Log-Likelihood Estimation*

$$\left| M(p(\mathbf{y})) - \hat{M}(p(\mathbf{y})) \right| \leq \epsilon \iff N \geq \frac{C_\epsilon}{P_\epsilon} \log \left(\frac{1}{\delta} \right), \quad (10)$$

where $\mathcal{Y}_\epsilon = \{ \mathbf{y} \in \mathcal{Y}_T \mid \log p(\mathbf{y}^*) - \log p(\mathbf{y}) \leq \epsilon \}$,
 $P_\epsilon = \sum_{\mathbf{y} \in \mathcal{Y}_\epsilon} p(\mathbf{y})$, $Q_\epsilon = \sum_{\mathbf{y} \in \mathcal{Y}_\epsilon} q(\mathbf{y})$, and $C_\epsilon = P_\epsilon / Q_\epsilon$.

2. *Sample-Complexity Bound for Importance-Weighted Entropy Estimation*

$$\left| H(p(\mathbf{y})) - \hat{H}(p(\mathbf{y})) \right| \leq \epsilon \iff N \geq \frac{(b-a)^2 C^2}{2\epsilon^2} \log \left(\frac{2}{\delta} \right), \quad (11)$$

where $p(\mathbf{y})/q(\mathbf{y}) \leq C$, $\forall \mathbf{y} \in \mathcal{Y}_T$.

Proof. 1: Probability of failure for N samples $(1 - Q_\epsilon)^N$, and $\log(1/(1-x)) \leq x^{-1}$ for $x \in (0, 1)$.
 2: Applying Hoeffding’s inequality after bounding importance weights by their worst-case value C . Detailed derivations are provided in Apx. B. \square

Discussion. The sample-complexity bound on estimating $M(p(\mathbf{y}))$ depends on the concentration of output sequences in the ϵ -region. This is desirable for autoregressive language generation, as sampling strategies generally focus on obtaining likely output sequences. In contrast, the sample complexity bound on estimating $H(p(\mathbf{y}))$ depends on the range of possible sequence likelihoods and the worst-case importance weight C squared. Both can be very high in practical settings. Thus, regarding sample-complexity, estimating $M(p(\mathbf{y}))$ appears desirable compared to estimating $H(p(\mathbf{y}))$ for typical LLM output distributions.

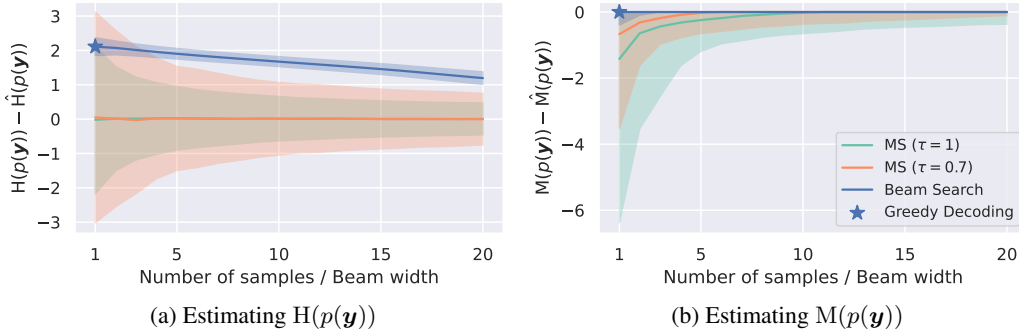


Figure 1: Quality of estimators for synthetic predictive distributions $p(\mathbf{y})$ with $|\mathcal{V}| = 20$ and $T = 4$. The predictive entropy $H(p(\mathbf{y}))$ is estimated as in Eq. (4) using multinomial sampling (MS) with different temperatures (τ). The negative maximum sequence log-likelihood $M(p(\mathbf{y}))$ is estimated by the maximum over N samples obtained by beam search ($N = 1$ represents greedy decoding) or MS with different τ . Statistics are obtained by sampling 2000 different $p(\mathbf{y})$. (a) Lines show average, shades denote one std. (b) Lines show the median, and shades denote the 5% to 95% quantile range.

Simulation Study. To gain insights into the practical implications of Theorem 1, we conduct a synthetic experiment where we sample predictive distributions $p(\mathbf{y})$ with smaller vocabulary sizes and shorter sequence lengths, while preserving the distributional characteristics typical for LLMs. This experiment design allows us to obtain ground truths for both $H(p(\mathbf{y}))$ and $M(p(\mathbf{y}))$, which become intractable to compute with larger vocabulary sizes and sequence lengths. Using this synthetic setup, we evaluate how the quality of estimators improves with the number of samples.

Fig. 1a summarizes the results for estimating $H(p(\mathbf{y}))$, derived from the logarithmic score. The results show that low sample sizes lead to high estimator variance. Similarly, Fig. 1b summarizes the results for estimating $M(p(\mathbf{y}))$, derived from the zero-one score. The results indicate that heuristics such as beam search and even greedy decoding provide accurate estimates of $M(p(\mathbf{y}))$ with high probability. In contrast, the variance of estimating $H(p(\mathbf{y}))$ remains substantial even when considering multiple samples and further increases when sampling with a lower temperature. Details on the sampling procedure of the predictive distributions, as well as additional experiments are provided in Apx. C.

4 RELATED WORK

In Sec. 1, we discussed prior work that considers single-sequence measures as heuristic baselines (Fadeeva et al., 2023; 2024; Vazhentsev et al., 2024; Abbasi-Yadkori et al., 2024; Plaut et al., 2024; Vashurin et al., 2025a;b), yet there is no clear consensus on how to apply them in a principled way. For instance, Ren et al. (2023) investigate length-normalized sequence likelihood (i.e., perplexity) and Zhang et al. (2025) use likelihood ratios between pre-trained and fine-tuned LLMs for OOD detection, but neither considers the MSP. Likewise, (Chen et al., 2024; Bakman et al., 2024; Yaldiz et al., 2025) include only perplexity as a baseline without specifying how the output sequence is obtained. Qiu and Miikkulainen (2024) also compare against a length-normalized variant, using the normalization from Murray and Chiang (2018), which corrects length bias for neural machine translation rather than uncertainty quantification. By grounding the MSP in theory and guiding its best approximation, we provide a foundation for advancing single-sequence uncertainty measures.

In Sec. 2.2, we also elaborated on the established sampling-based uncertainty estimation measures, i.e., PE (Malinin and Gales, 2018) and SE (Kuhn et al., 2023; Farquhar et al., 2024). There is a body of work that builds upon the concept of PE, for instance, by considering a weighting factor for individual token and sequence likelihoods to account for the importance on a semantic level (Duan et al., 2024; Bakman et al., 2024; Yaldiz et al., 2025). There is also a body of work that extends the concept of SE (Kuhn et al., 2023; Farquhar et al., 2024), for instance, by improving the semantic clustering (Nikitin et al., 2024; Qiu and Miikkulainen, 2024), improving the sampling of diverse output sequences (Aichberger et al., 2025), or directly approximating the measure from hidden states of the language model (Kossen et al., 2024; Chen et al., 2024). Since our goal is to compare the uncertainty principles induced by the underlying proper scoring rules, we focus on their established forms for a fair and consistent comparison.

There is also work on uncertainty estimation in NLG that cannot be directly grounded in proper scoring rules. For instance, several approaches leverage the language model itself to directly predict uncertainty, either through numerical estimates or verbal explanations (Mielke et al., 2022; Lin et al., 2022; Kadavath et al., 2022; Cohen et al., 2023a; Ganguli et al., 2023; Ren et al., 2023; Tian et al., 2023). Cohen et al. (2023b) use cross-examination, where one model generates an output, and another evaluates it, while Manakul et al. (2023) assess uncertainty by feeding multiple sampled outputs into a second model. Zhou et al. (2023) explore the behavior of LLMs when expressing their uncertainty, providing insights into how models articulate confidence. Xiao et al. (2022) analyze how factors such as model architecture and training data affect uncertainty estimates. Gao et al. (2024) perturb inputs to quantify uncertainty, and Chen and Mueller (2024) rely on consistency and self-reflection under different prompting strategies. Jiang et al. (2024) propose Graph Uncertainty, viewing individual claims in an output sequence as nodes in a bipartite graph. Finally, conformal prediction (Quach et al., 2024) provides a calibration-based rule for determining when to stop sampling during generation.

5 EXPERIMENTS

We aligned our experiments on evaluating uncertainty estimation methods with prior work by focusing on free-form question answering tasks (Kuhn et al., 2023; Duan et al., 2024; Bakman et al., 2024; Nikitin et al., 2024; Aichberger et al., 2025; Kossen et al., 2024). While Farquhar et al. (2024) additionally concerns experiments with paragraph-length generations, their approach breaks down the paragraph into factual claims and constructs corresponding questions. Therefore, performance on this task is expected to align with general free-form question answering tasks, and we therefore focused on those for a clearer and more direct evaluation. This focus further avoids potential confounding factors introduced by additional experimental complexities.

Datasets. We evaluated uncertainty estimation methods on three different datasets. We used the more than 3,000 test instances from *TriviaQA* (Joshi et al., 2017) concerning trivia questions, the more than 300 test instances from *SVAMP* (Patel et al., 2021) concerning elementary-level math problems, and the more than 3,600 test instances from *NQ-Open* (Lee et al., 2019) to assess natural questions aggregated from Google Search. Each dataset was used for two distinct tasks: (1) generating concise answers in the form of short phrases (*short*) and (2) generating more detailed answers in the form of full sentences (*long*), following the experimental setup in Farquhar et al. (2024). The six individual tasks were selected so that they align with prior work to provide a fair and meaningful comparison. They aim at covering a broad range of real-world scenarios, ensuring a comprehensive evaluation.

Models. We conducted our evaluations on six different LLMs covering various architectures, sizes, and training stages. Specifically, we used the transformer model series *Llama-3.1* (Dubey et al., 2024) and the state-space model series *Falcon Mamba* (Gu and Dao, 2024; Zuo et al., 2024), representing two prominent paradigms. To assess the effect of training stage and scale on uncertainty estimation in NLG, we considered pre-trained (*PT*) and instruction-tuned (*IT*) LLMs with 7, 8, and 70 billion parameters, covering a wide spectrum of model characteristics used in real-world applications.

Baselines. We compare our method, G-NLL, against the commonly used uncertainty measures based on the logarithmic score as of Eq. (4) and Eq. (5) and their variants. These include predictive entropy (*PE*), length-normalized predictive entropy (*LN-PE*) (Malinin and Gales, 2021), semantic entropy (*SE*), length-normalized semantic entropy (*LN-SE*), and Discrete semantic entropy (*D-SE*) (Kuhn et al., 2023; Farquhar et al., 2024). For a given output sequence \mathbf{y}' , the length-normalized variants consider $\bar{p}(\mathbf{y}' | \mathbf{x}, \mathbf{w})$ instead of $p(\mathbf{y}' | \mathbf{x}, \mathbf{w})$ to compute the uncertainty estimates. D-SE completely disregards the likelihood of the output sequence and only considers the proportion of output sequences that belong to the same semantic cluster (Farquhar et al., 2024). These baselines represent direct counterparts derived from alternative proper scoring rules as discussed in Sec. 2.2. More importantly, they are among the most widely used uncertainty estimation methods in NLG. Comparing against these established methods provides a decisive proxy for the empirical performance of G-NLL.

Evaluation. Effective uncertainty measures should accurately reflect the reliability of answers generated by the language model. In other words, higher uncertainty should correspond to a greater tendency for incorrect outputs. Thus, to evaluate the performance of an uncertainty estimator, we assess how well the uncertainty estimate correlates with the correctness of the language model’s answer. Correct answers should be assigned a lower uncertainty estimator than incorrect answers. We consider a given answer correct if the F1 score of the commonly used SQuAD metric to the ground

Table 1: Average AUROC (\uparrow) with standard errors across TriviaQA, SVAMP, and NQ datasets, using uncertainty estimates to distinguish between correct and incorrect answers. Six different language models, with varying model architectures (*transformer*, *state-space*), model sizes (*7B*, *8B*, *70B*), and training stages (*PT*, *IT*) are considered. The reference answer is generated using greedy decoding, either as a whole sentence (*long*) or a short phrase (*short*). Its correctness is assessed by *F1* score using SQuAD > 0.5 as decision rule or LLM-as-a-judge (*LLM*). *PE*, *LN-PE*, *SE*, *LN-SE*, and *D-SE* use 10 output sequences (generated via multinomial sampling) to obtain their uncertainty estimates. G-NLL solely uses one output sequence (generated via greedy decoding) for its uncertainty estimate. * indicates that G-NLL performs significantly better than the best sampling-based method ($p < 0.05$).

Underlying proper scoring rule			Logarithmic					Zero-One		
Language Model	Gen.	Metric	PE	LN-PE	SE	LN-SE	D-SE	G-NLL		
Transformer	8B	PT	short F1	.776 \pm .009	.795 \pm .008	.775 \pm .008	.793 \pm .008	<u>.804</u> \pm .009	.824 \pm .008*	
		PT	short LLM	.698 \pm .014	.714 \pm .014	.690 \pm .014	.706 \pm .014	<u>.719</u> \pm .014	.726 \pm .014	
		PT	long LLM	.562 \pm .012	.555 \pm .012	.545 \pm .013	.553 \pm .012	<u>.600</u> \pm .012	.649 \pm .012*	
	8B	IT	short F1	.772 \pm .008	.801 \pm .007	.805 \pm .007	.814 \pm .007	.806 \pm .008	.838 \pm .006*	
		IT	short LLM	.676 \pm .015	.697 \pm .015	.704 \pm .015	<u>.709</u> \pm .015	.694 \pm .016	.722 \pm .014	
		IT	long LLM	.551 \pm .012	.548 \pm .012	.599 \pm .012	<u>.601</u> \pm .012	<u>.609</u> \pm .012	.615 \pm .012	
	70B	PT	short F1	.775 \pm .010	.790 \pm .010	.793 \pm .009	<u>.803</u> \pm .009	.791 \pm .009	.820 \pm .008	
			PT	short LLM	.693 \pm .019	.709 \pm .019	.718 \pm .019	<u>.722</u> \pm .020	.715 \pm .018	.723 \pm .019
			PT	long LLM	.552 \pm .014	.534 \pm .014	.558 \pm .015	.569 \pm .013	<u>.571</u> \pm .013	.649 \pm .012*
		IT	short F1	.748 \pm .010	.781 \pm .009	.790 \pm .009	.799 \pm .009	.783 \pm .009	.792 \pm .011	
			IT	short LLM	.681 \pm .021	.698 \pm .022	.703 \pm .022	.709 \pm .022	.699 \pm .019	.699 \pm .024
			IT	long LLM	.555 \pm .013	.557 \pm .014	.568 \pm .014	.595 \pm .014	.600 \pm .014	.562 \pm .014
State-Space	7B	PT	short F1	.811 \pm .007	.815 \pm .007	.809 \pm .008	.822 \pm .008	<u>.828</u> \pm .006	.843 \pm .006*	
		PT	short LLM	.705 \pm .012	.711 \pm .011	.701 \pm .012	.711 \pm .012	<u>.716</u> \pm .012	.728 \pm .011	
		PT	long LLM	.567 \pm .012	.597 \pm .012	.574 \pm .012	.611 \pm .012	.624 \pm .012	.612 \pm .012	
	IT	short F1	.793 \pm .009	.814 \pm .007	.797 \pm .008	.816 \pm .008	<u>.829</u> \pm .007	.838 \pm .006		
		IT	short LLM	.690 \pm .013	.701 \pm .012	.689 \pm .012	.699 \pm .012	<u>.711</u> \pm .013	.719 \pm .012	
		IT	long LLM	.588 \pm .012	.587 \pm .012	.597 \pm .012	.618 \pm .012	.629 \pm .012	.615 \pm .012	
Average			.677 \pm .003	.689 \pm .003	.690 \pm .003	.703 \pm .003	<u>.707</u> \pm .003	.721 \pm .003*		

truth answer exceeds 0.5 (Rajpurkar et al., 2016). Since this metric is only applicable for short-phrase generations that align with the ground truth answer, we additionally employ Llama-3.1 with 70 billion parameters (Dubey et al., 2024) as LLM-as-a-judge to assess the correctness of both short-phrase and full-sentence generations. Finally, to measure the correlation between the incorrectness of answers and the respective uncertainty estimates, we use the AUROC. Higher AUROC values indicate better performance of the uncertainty estimator, as it reflects a stronger alignment between the correctness of the language model’s answers and their respective uncertainty estimates. This evaluation process follows established protocols for uncertainty estimation in NLG (Kuhn et al., 2023; Farquhar et al., 2024; Duan et al., 2024; Bakman et al., 2024; Nikitin et al., 2024; Kossen et al., 2024).

5.1 MAIN RESULTS

Tab. 1 summarizes the performance of the uncertainty measures across the six LLMs, six tasks, and two evaluation metrics. We report the average AUROC across the three datasets, highlighting the best-performing measure in bold. Additionally, the best-performing measure based on the logarithmic score is underlined, unless it also represents the overall best. In 13 out of 18 scenarios, G-NLL outperforms the uncertainty measures and remains competitive in the remaining 5 scenarios. This strong performance is particularly evident in tasks involving the generation of short phrases, suggesting its effectiveness in capturing the essential part of the output sequence that contains the factual answer to a question. This is especially valuable in practical scenarios, where the uncertainty about a specific fact is often more critical than the uncertainty about the entire generated sentence. Overall, our measure significantly outperforms all other measures when considering the average across all scenarios. This demonstrates that G-NLL achieves strong empirical performance despite relying on only a single output sequence. Crucially, G-NLL relies only on the greedily decoded output sequence, yielding a deterministic and hyperparameter-free uncertainty estimate. Detailed evaluation results are provided in Apx. D.

5.2 ABLATION STUDY ON APPROXIMATION VARIANTS

The strong empirical performance of G-NLL on question answering NLG tasks suggests that it effectively captures key aspects of uncertainty, even when using significantly fewer output sequences compared to uncertainty measures based on the logarithmic score. To examine the underlying factors driving this behavior, we analyze the impact of the sampling method on approximating the measure of aleatoric uncertainty $M(p(\mathbf{y} | \mathbf{x}, \mathbf{w}))$ derived from the zero-one score in Eq. (7).

Specifically, we utilize multinomial sampling (MS) with different temperatures τ , beam search (BS) with different beam sizes, and greedy decoding as for G-NLL. For each sampling method, we generate a single output sequence per instance in the TriviaQA dataset. We then use the corresponding NLL as the uncertainty estimate, following the same evaluation process as in the main experiments above. Notably, the baselines are unaffected by the choice of sampling method for computing the NLL, as we again use their optimal hyperparameter settings for sampling.

The results in Fig. 2 show that better approximations of the most likely output sequences indeed lead to higher uncertainty estimation performance, reinforcing the validity of our alternative measure derived from the zero-one score. Additionally, it can be observed that sampling output sequences using greedy decoding significantly outperforms MS. While performance improves further with BS, as anticipated, the marginal benefits are relatively small. This supports the claim that greedy decoding provides a strong approximation to the most likely output sequence. Since BS is computationally more expensive (as its beam size corresponds to the number of sampled output sequences), using greedy decoding in G-NLL achieves the best trade-off between effectiveness and efficiency.

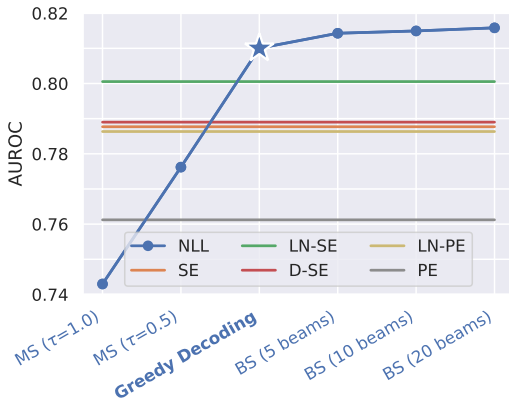


Figure 2: Average AUROC (\uparrow) using uncertainty measures based on the zero-one score, with the output sequence to approximate the MSP being generated with multinomial sampling (MS), greedy decoding (G-NLL), or beam search (BS), compared to baseline measures based on the logarithmic score.

6 CONCLUSION

In this work, we theoretically motivate an alternative uncertainty measure: the NLL of the most likely output sequence under a given language model. The measure is grounded in the general notion of proper scoring rules, with the zero-one score serving as an alternative to the commonly adopted logarithmic score. Unlike multi-sequence uncertainty measures, it can be efficiently approximated with G-NLL from a single greedily decoded sequence, challenging the widespread use of sampling- and clustering-based approaches for uncertainty estimation in NLG. Experiments show that G-NLL outperforms prior methods that entail considerably higher computational costs and algorithmic complexity. Importantly, we find that the choice of decoding strategy to approximate the most likely output sequence is critical, and sampling-based or length-normalized variants can degrade approximation quality.

While G-NLL effectively captures uncertainty, it does not explicitly account for semantics. Exploring extensions that incorporate semantic information could further enhance single-sequence uncertainty estimation, but would likely increase computational cost (see Sec. 2.3), and we leave exploring this trade-off for future work. Moreover, G-NLL requires access to the probabilities of generated tokens, though not the full token distribution. While such information is typically exposed by modern LLM APIs, future work may investigate how to approximate it when unavailable.

While there remain opportunities for refinement, G-NLL provides a principled foundation for efficient uncertainty estimation in NLG. Its simplicity and minimal computational overhead make it a practical baseline for developing new uncertainty measures and support scalable deployment in real-world applications. More broadly, our work shows that single-sequence measures serve as a viable alternative to sampling-based methods, offering a principled way to estimate uncertainty efficiently in LLMs.

ACKNOWLEDGEMENTS

We acknowledge EuroHPC Joint Undertaking for awarding us access to Karolina at IT4Innovations, Czech Republic. The ELLIS Unit Linz, the LIT AI Lab, the Institute for Machine Learning, are supported by the Federal State Upper Austria. We thank the projects FWF AIRI FG 9-N (10.55776/FG9), AI4GreenHeatingGrids (FFG- 899943), Stars4Waters (HORIZON-CL6-2021-CLIMATE-01-01), FWF Bilateral Artificial Intelligence (10.55776/COE12). We thank NXAI GmbH, Audi AG, Silicon Austria Labs (SAL), Merck Healthcare KGaA, GLS (Univ. Waterloo), TÜV Holding GmbH, Software Competence Center Hagenberg GmbH, dSPACE GmbH, TRUMPF SE + Co. KG.

ETHICS STATEMENT

This work contributes to advancing the efficiency and reliability of uncertainty estimation in LLMs. Our proposed measure, G-NLL, reduces computational overhead compared to existing approaches, enabling broader adoption of uncertainty estimation techniques in real-world applications. While we recognize societal and ethical risks associated with LLMs, such as misinformation, our contribution seeks to mitigate these risks by supporting more trustworthy deployment through improved uncertainty estimation. Our work does not involve human subjects, sensitive data, or personally identifiable information. All datasets used are publicly available and commonly employed in prior work on NLG.

REPRODUCIBILITY STATEMENT

We provide detailed descriptions of datasets, experimental setup, and evaluation metrics in Sec. 5 of the main paper and Apx. D. All datasets used are publicly available, and we utilize standard benchmarks to support straightforward replication. All experiments were performed on a single node with 8 NVIDIA A100 Tensor Core GPUs. Running all evaluations required roughly 100 node hours.

REFERENCES

- Yasin Abbasi-Yadkori, Ilja Kuzborskij, András György, and Csaba Szepesvari. To believe or not to believe your LLM: Iterative prompting for estimating epistemic uncertainty. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. Improving uncertainty estimation through semantically diverse language generation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Yavuz Faruk Bakman, Duygu Nur Yaldiz, Baturalp Buyukates, Chenyang Tao, Dimitrios Dimitriadis, and Salman Avestimehr. MARS: Meaning-aware response scoring for uncertainty estimation in generative LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pages 7752–7767. Association for Computational Linguistics, 2024.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. INSIDE: LLMs’ internal states retain the power of hallucination detection. In *The Twelfth International Conference on Learning Representations*, 2024.
- Jiuhai Chen and Jonas Mueller. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2024.
- Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. Crawling the internal knowledge-base of language models. In Andreas Vlachos and Isabelle Augenstein, editors, *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1856–1869. Association for Computational Linguistics, 2023a.
- Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. LM vs LM: Detecting factual errors via cross examination. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the*

- 2023 *Conference on Empirical Methods in Natural Language Processing*, pages 12621–12640. Association for Computational Linguistics, 2023b.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pages 5050–5063, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and Aurelien Rodriguez et al. The llama 3 herd of models. *arXiv*, 2407.21783, 2024.
- 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. LM-polygraph: Uncertainty estimation for language models. In Yansong Feng and Els Lefever, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 446–461. Association for Computational Linguistics, 2023.
- Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. Fact-checking the output of large language models via token-level uncertainty quantification. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9367–9385. Association for Computational Linguistics, 2024.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
- Yarin Gal. *Uncertainty in deep learning*. PhD thesis, University of Cambridge, 2016.
- Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I. Liao, Kamilė Lukošiuūtė, Anna Chen, Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr, Jared Mueller, Joshua Landau, Kamal Ndousse, Karina Nguyen, Liane Lovitt, Michael Sellitto, Nelson Elhage, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robert Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu, Saurav Kadavath, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Christopher Olah, Jack Clark, Samuel R. Bowman, and Jared Kaplan. The capacity for moral self-correction in large language models. *arXiv*, 2302.07459, 2023.
- Xiang Gao, Jiaxin Zhang, Lalla Mouatadid, and Kamalika Das. SPUQ: Perturbation-based uncertainty quantification for large language models. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2024.
- Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.
- Sebastian Gruber and Florian Buettner. Uncertainty estimates of predictions via a general bias-variance decomposition. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206, pages 11331–11354. PMLR, 2023.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*, 2024.
- Paul Hofman, Yusuf Sale, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty with proper scoring rules. *arXiv*, 2404.12215, 2024.

- Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv*, 1112.5745, 2011.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.
- Mingjian Jiang, Yangjun Ruan, Prasanna Sattigeri, Salim Roukos, and Tatsunori Hashimoto. Graph-based uncertainty metrics for long-form language model generations. In *Advances in Neural Information Processing Systems*, 2024.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2017.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. Language models (mostly) know what they know. *arXiv*, 2207.05221, 2022.
- Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa Schut, Shreshth Malik, and Yarin Gal. Semantic entropy probes: Robust and cheap hallucination detection in llms. *arXiv*, 2406.15927, 2024.
- Nikita Kotelevskii, Vladimir Kondratyev, Martin Takáč, Eric Moulines, and Maxim Panov. From risk to uncertainty: Generating predictive uncertainty measures via bayesian estimation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*, 2023.
- Meelis Kull and Peter Flach. Novel decompositions of proper scoring rules for classification: Score adjustment as precursor to calibration. In Annalisa Appice, Pedro Pereira Rodrigues, Vítor Santos Costa, Carlos Soares, João Gama, and Alípio Jorge, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 68–85. Springer International Publishing, 2015.
- Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym Korablyov, and Yoshua Bengio. DEUP: Direct epistemic uncertainty prediction. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096. Association for Computational Linguistics, 2019.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856.
- Andrey Malinin. *Uncertainty estimation in deep learning with application to spoken language assessment*. PhD thesis, University of Cambridge, 2019.
- Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In *International Conference on Learning Representations*, 2021.

- Potsawee Manakul, Adian Liusie, and Mark Gales. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017. Association for Computational Linguistics, 2023.
- Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational agents’ overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872, 2022.
- Kenton Murray and David Chiang. Correcting length bias in neural machine translation. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéal, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Research Papers*, 2018.
- Alexander V Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. Kernel language entropy: Fine-grained uncertainty quantification for LLMs from semantic similarities. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094. Association for Computational Linguistics, 2021.
- Benjamin Plaut, Nguyen X. Khanh, and Tu Trinh. Probabilities of chat llms are miscalibrated but still predict correctness on multiple-choice q&a. *arXiv*, 2402.13213, 2024.
- Xin Qiu and Risto Miikkulainen. Semantic density: Uncertainty quantification for large language models through confidence measurement in semantic space. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Victor Quach, Adam Fisch, Tal Schuster, Adam Yala, Jae Ho Sohn, Tommi S. Jaakkola, and Regina Barzilay. Conformal language modeling. In *The Twelfth International Conference on Learning Representations*, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2018.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392. Association for Computational Linguistics, 2016.
- Jie Ren, Yao Zhao, Tu Vu, Peter J. Liu, and Balaji Lakshminarayanan. Self-evaluation improves selective generation in large language models. *arXiv*, 2312.09300, 2023.
- Alfréd Rényi. On measures of entropy and information. In *Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics and Probability*, pages 547–561. University of California Press, 1961.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. Introducing an improved information-theoretic measure of predictive uncertainty. *arXiv*, 2311.08309, 2023a.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, Günter Klambauer, and Sepp Hochreiter. Quantification of uncertainty with adversarial models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 19446–19484. Curran Associates, Inc., 2023b.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. On information-theoretic measures of predictive uncertainty. *Uncertainty in Artificial Intelligence*, 2025.

- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442. Association for Computational Linguistics, 2023.
- Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Daniil Vasilev, Akim Tsvigun, Sergey Petrakov, Rui Xing, Abdelrahman Sadallah, Kirill Grishchenkov, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, Maxim Panov, and Artem Shelmanov. Benchmarking uncertainty quantification methods for large language models with LM-polygraph. *Transactions of the Association for Computational Linguistics*, 2025a.
- Roman Vashurin, Maiya Goloburda, Preslav Nakov, Artem Shelmanov, and Maxim Panov. Cocoa: A generalized approach to uncertainty quantification by integrating confidence and consistency of llm outputs. *arXiv*, 2502.04964, 2025b.
- Artem Vazhentsev, Ekaterina Fadeeva, Rui Xing, Alexander Panchenko, Preslav Nakov, Timothy Baldwin, Maxim Panov, and Artem Shelmanov. Unconditional truthfulness: Learning conditional dependency for uncertainty quantification of large language models. *arXiv*, 2408.10692, 2024.
- Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scale empirical analysis. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7273–7284. Association for Computational Linguistics, 2022.
- Duygu Nur Yaldiz, Yavuz Faruk Bakman, Baturalp Buyukates, Chenyang Tao, Anil Ramakrishna, Dimitrios Dimitriadis, Jieyu Zhao, and Salman Avestimehr. Do not design, learn: A trainable scoring function for uncertainty estimation in generative LLMs. In *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025.
- Andi Zhang, Tim Z. Xiao, Weiyang Liu, Robert Bamler, and Damon Wischik. Your finetuned large language model is already a powerful out-of-distribution detector. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*. PMLR, 2025.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models. *arXiv*, 2205.01068, 2022.
- Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. Navigating the grey area: How expressions of uncertainty and overconfidence affect language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5506–5524. Association for Computational Linguistics, 2023.
- Jingwei Zuo, Maksim Velikanov, Dhia Eddine Rhaiem, Ilyas Chahed, Younes Belkada, Guillaume Kunsch, and Hakim Hacid. Falcon mamba: The first competitive attention-free 7b language model. *arXiv*, 2410.05355, 2024.

USAGE OF LARGE LANGUAGE MODELS (LLMs)

LLMs were used as general-purpose assistive tools during the preparation of this paper. Their use fell into two categories: (i) writing assistance, where they were used to improve the clarity and readability of certain passages through language refinement, and (ii) coding assistance, where they were used for support with code completion and debugging. LLMs were not used for research ideation, experimental design, theoretical development, or result analysis. All substantive contributions, including the conception of ideas, methodology, and experiments, were made by the authors.

A DETAILS ON PREDICTIVE UNCERTAINTY IN NLG

As stated in Sec. 2, the total uncertainty of a given language model is defined as

$$\mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}'))]] .$$

This corresponds to the expected score that models from the posterior $p(\tilde{\mathbf{w}} | \mathcal{D})$ assign to output sequences \mathbf{y}' drawn from the predictive distribution of the given model $p(\mathbf{y}' | \mathbf{x}, \mathbf{w})$. Note again that \mathbf{y}' is an independent notational copy of \mathbf{y} introduced to make the random variable to integrate within the expectation more explicit.

Adding and subtracting $\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}'))]$ yields the general scoring-rule-based uncertainty decomposition in Eq. (1) of the main paper:

$$\begin{aligned} & \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}'))]]}_{\text{expected score}} = \\ & \mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}'))] \\ & \quad + \mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}'))]] - \mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}'))] = \\ & \underbrace{\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}'))]}_{\text{entropy term}} + \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [\mathbb{S}(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}')) - \mathbb{S}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}'))]}_{\text{divergence term}} . \end{aligned}$$

A.1 UNCERTAINTY MEASURES BASED ON THE LOGARITHMIC SCORE

To derive the established uncertainty measures in Sec. 2.2, we substitute the logarithmic score

$$S_{\log}(p(\mathbf{y} | \mathbf{x}, \cdot), \mathbf{y}') = -\log p(\mathbf{y} = \mathbf{y}' | \mathbf{x}, \cdot) \quad (2)$$

into the general scoring-rule-based uncertainty decomposition (Eq. (1)) to get:

$$\begin{aligned} & \mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \tilde{\mathbf{w}})]] = \\ & \mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \mathbf{w})] + \mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \tilde{\mathbf{w}}) + \log p(\mathbf{y}' | \mathbf{x}, \mathbf{w})]] . \end{aligned}$$

On the LHS, we note the definition of the cross-entropy

$$\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \tilde{\mathbf{w}})] = \text{CE}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) ; p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}})) .$$

On the RHS, we note that the first term is the definition of the Shannon entropy

$$\mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \mathbf{w})] = \text{H}(p(\mathbf{y} | \mathbf{x}, \mathbf{w})) ,$$

and the second term is the definition of the Kullback-Leibler divergence

$$\begin{aligned} \mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} [-\log p(\mathbf{y}' | \mathbf{x}, \tilde{\mathbf{w}}) + \log p(\mathbf{y}' | \mathbf{x}, \mathbf{w})] &= \mathbb{E}_{p(\mathbf{y}'|\mathbf{x},\mathbf{w})} \left[\frac{\log p(\mathbf{y}' | \mathbf{x}, \mathbf{w})}{\log p(\mathbf{y}' | \mathbf{x}, \tilde{\mathbf{w}})} \right] \\ &= \text{KL}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) \| p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}})) . \end{aligned}$$

Using these definitions directly results in Eq. (3) of the main paper:

$$\begin{aligned} & \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\text{CE}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) ; p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{total uncertainty}} = \quad (3) \\ & \underbrace{\text{H}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}))}_{\text{aleatoric uncertainty}} + \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}}|\mathcal{D})} [\text{KL}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) \| p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{epistemic uncertainty}} . \end{aligned}$$

A.2 UNCERTAINTY MEASURES BASED ON THE ZERO-ONE SCORE

To derive the new uncertainty measures in Sec. 2.3, we substitute the zero-one score

$$S_{0-1}(p(\mathbf{y} | \mathbf{x}, \cdot), \mathbf{y}') = 1 - \mathbb{1}\{\mathbf{y}' = \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{y} | \mathbf{x}, \cdot)\} \quad (6)$$

into the general scoring-rule-based uncertainty decomposition (Eq. (1)).

For clarity, we distinguish between the most likely output sequence under the given model \mathbf{w}

$$\mathbf{y}^* := \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{y} | \mathbf{x}, \mathbf{w}),$$

and the most likely sequence under a model $\tilde{\mathbf{w}}$ drawn from the posterior $p(\tilde{\mathbf{w}} | \mathcal{D})$

$$\tilde{\mathbf{y}}^* := \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}).$$

Because the indicator enforces $\mathbf{y}' = \mathbf{y}^*$ or $\mathbf{y}' = \tilde{\mathbf{y}}^*$, the inner expectation over \mathbf{y}' collapses, yielding

$$\begin{aligned} \mathbb{E}_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})} [S_{0-1}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}), \mathbf{y}')] &= 1 - p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w}), \\ \mathbb{E}_{p(\mathbf{y}' | \mathbf{x}, \mathbf{w})} [S_{0-1}(p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}), \mathbf{y}')] &= 1 - p(\mathbf{y} = \tilde{\mathbf{y}}^* | \mathbf{x}, \mathbf{w}). \end{aligned}$$

Using these definitions directly results in Eq. (7) of the main paper:

$$\begin{aligned} \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}} | \mathcal{D})} [1 - p(\mathbf{y} = \tilde{\mathbf{y}}^* | \mathbf{x}, \mathbf{w})]}_{\text{total uncertainty}} &= \quad (7) \\ \underbrace{1 - p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w})}_{\text{aleatoric uncertainty}} &+ \underbrace{p(\mathbf{y} = \mathbf{y}^* | \mathbf{x}, \mathbf{w}) - \mathbb{E}_{p(\tilde{\mathbf{w}} | \mathcal{D})} [p(\mathbf{y} = \tilde{\mathbf{y}}^* | \mathbf{x}, \mathbf{w})]}_{\text{epistemic uncertainty}}. \end{aligned}$$

Incorporating Semantics. As discussed in Sec. 2.1, uncertainty measures are determined by a predictive distribution together with a proper scoring rule. In place of the output-sequence distribution $p(\mathbf{y} | \mathbf{x}, \mathbf{w})$, we may therefore consider the semantic cluster distribution $p(c | \mathbf{x}, \mathbf{w})$ to incorporate semantic information into the uncertainty measures based on the zero-one score.

In this setting, the zero-one score evaluates the predictive distribution at its *most likely semantic cluster* rather than at the most likely output sequence:

$$S_{0-1}(p(c | \mathbf{x}, \cdot), c') = 1 - \mathbb{1}\{c' = \underset{c}{\operatorname{argmax}} p(c | \mathbf{x}, \cdot)\}. \quad (12)$$

Again, c' is an independent notational copy of c introduced to make clear which variable the expectation is calculated for. For clarity, we again distinguish between the most likely semantic cluster under the given model and that under a model drawn from the posterior:

$$c^* := \underset{c}{\operatorname{argmax}} p(c | \mathbf{x}, \mathbf{w}), \quad \tilde{c}^* := \underset{c}{\operatorname{argmax}} p(c | \mathbf{x}, \tilde{\mathbf{w}}).$$

Substituting the zero-one score defined in Eq. (12) into the general scoring-rule-based uncertainty decomposition (Eq. (1)) yields the following decomposition of total semantic uncertainty, the expected confidence that the given model assigns to the most likely *semantic cluster* predicted by models drawn from the posterior $p(\tilde{\mathbf{w}} | \mathcal{D})$:

$$\begin{aligned} \underbrace{\mathbb{E}_{p(\tilde{\mathbf{w}} | \mathcal{D})} [1 - p(c = \tilde{c}^* | \mathbf{x}, \mathbf{w})]}_{\text{total semantic uncertainty}} &= \quad (13) \\ \underbrace{1 - p(c = c^* | \mathbf{x}, \mathbf{w})}_{\text{aleatoric semantic uncertainty}} &+ \underbrace{p(c = c^* | \mathbf{x}, \mathbf{w}) - \mathbb{E}_{p(\tilde{\mathbf{w}} | \mathcal{D})} [p(c = \tilde{c}^* | \mathbf{x}, \mathbf{w})]}_{\text{epistemic semantic uncertainty}}. \end{aligned}$$

As in Eq. (7), the epistemic semantic uncertainty is a posterior expectation that is challenging to estimate. The aleatoric semantic uncertainty considers the likelihood of the most likely *semantic cluster* under the given language model, i.e, the most-likely cluster probability (MCP), analogous to maximum sequence probability (MSP).

However, approximating the most likely semantic cluster is more challenging, since the semantic cluster distribution $p(c | \mathbf{x}, \mathbf{w})$ cannot be sampled from directly (Aichberger et al., 2025). Whether MCP admits an efficient approximation remains an open question. While MSP (as approximated by G-NLL) may correlate with MCP, a detailed analysis of this relationship is left to future work.

B DETAILS ON THEORETICAL ANALYSIS

We want to compare the sample complexity of approximating the maximum and expected log-likelihood, showing that approximating the maximum is desirable in the setting of LLMs. In the following, we derive probabilistic concentration bounds on the number of samples N needed to approximate each quantity to within ϵ -precision with high probability and use these to reason about the sample complexity required for reliable approximation.

Setup. As detailed in Sec. 2 in the main paper, we consider the probability distribution $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ of output sequences \mathbf{y} induced by the LLM with parameters \mathbf{w} for an input sequence \mathbf{x} . We abbreviate this distribution by $p(\mathbf{y})$. We want to study the approximation error for the negative maximum log-likelihood $M(p(\mathbf{y})) = -\max_{\mathbf{y}} \log p(\mathbf{y})$ (min-entropy) and the negative expected log-likelihood $H(p(\mathbf{y})) = -\mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y})]$ (Shannon entropy) using N samples.

We assume boundedness, i.e. there exist constants a, b such that $a \leq \log p(\mathbf{y}) \leq b$, $\forall \mathbf{y} \in \mathcal{Y}_T$. This assumption generally holds in practice as softmax logits are typically clipped and the softmax temperature is non-zero, though the constants can be of large magnitude depending on the sequence length T . Furthermore, to incorporate importance sampling schemes (e.g., low-temperature sampling, top-k or top-p sampling, beam search, etc.), we consider access to a proposal distribution $q(\mathbf{y})$. This helps us understand how practical decoding choices impact concentration behavior.

B.1 APPROXIMATING THE MAXIMUM LOG-LIKELIHOOD (MIN-ENTROPY)

Let $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y})$ and $P_\epsilon = \sum_{\mathbf{y} \in \mathcal{Y}_\epsilon} p(\mathbf{y})$ as well as $Q_\epsilon = \sum_{\mathbf{y} \in \mathcal{Y}_\epsilon} q(\mathbf{y})$ where $\mathcal{Y}_\epsilon = \{\mathbf{y} \in \mathcal{Y}_T \mid \log p(\mathbf{y}^*) - \log p(\mathbf{y}) \leq \epsilon\}$. Thus \mathcal{Y}_ϵ is the set of output sequences with a log-likelihood within ϵ of the maximum log-probability, and P_ϵ, Q_ϵ are the cumulative probability masses under sampling distributions p and q of those output sequences. The correction factor $C_\epsilon = P_\epsilon/Q_\epsilon$ compares the ϵ -regions under distributions p and q . If they match, $C_\epsilon \approx 1$, if q is less likely to provide samples within the ϵ -region, $C_\epsilon > 1$ and if it is more likely $C_\epsilon < 1$. We empirically estimate $M(p(\mathbf{y}))$ with $\hat{M}(p(\mathbf{y})) = -\max\{\log p(\mathbf{y}^n)\}_{n=1}^N$, the maximum over the log-probabilities of the sampled output sequences.

When sampling from $q(\mathbf{y})$, the probability that no $\mathbf{y} \in \mathcal{Y}_\epsilon$ is not obtained in N samples is

$$(1 - Q_\epsilon)^N \leq \delta. \quad (14)$$

Solving for N , we obtain

$$N \geq \frac{\log(1/\delta)}{\log(1/(1 - Q_\epsilon))}. \quad (15)$$

Furthermore, we use $Q_\epsilon = P_\epsilon/C_\epsilon$ and the inequality $\log(1/(1-x)) \leq x^{-1}$ for $x \in (0, 1)$ to simplify further.

Result. To ensure with probability at least $1 - \delta$ that $\left| M(p(\mathbf{y})) - \hat{M}(p(\mathbf{y})) \right| \leq \epsilon$ it suffices that

$$N \geq \frac{C_\epsilon}{P_\epsilon} \log \left(\frac{1}{\delta} \right). \quad (16)$$

We will discuss the practical implications of this bound in Sec. B.3.

B.2 APPROXIMATING THE EXPECTED LOG-LIKELIHOOD (SHANNON ENTROPY)

We consider importance sampling with $q(\mathbf{y})$, thus

$$H(p(\mathbf{y})) = -\mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y})] = \mathbb{E}_{q(\mathbf{y})} \left[\frac{p(\mathbf{y})}{q(\mathbf{y})} \log p(\mathbf{y}) \right]$$

and its MC estimator

$$\hat{H}(p(\mathbf{y})) = -\frac{1}{N} \sum_{n=1}^N \frac{p(\mathbf{y}^n)}{q(\mathbf{y}^n)} \log p(\mathbf{y}^n)$$

with $\mathbf{y}^n \sim q(\mathbf{y})$. Furthermore, let $c(\mathbf{y}) = \frac{p(\mathbf{y})}{q(\mathbf{y})}$ be bounded, i.e. $c(\mathbf{y}) \leq C, \forall \mathbf{y} \in \mathcal{Y}_T$.

Hoeffding’s inequality states that

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{n=1}^N X_n - \mathbb{E}[X_n] \right| \geq \epsilon \right) \leq 2 \exp \left(-\frac{2N\epsilon^2}{(b-a)^2} \right), \quad (17)$$

for a random variable X_n with $a \leq X_n \leq b$. We set $X_n = c(\mathbf{y}^n) \log p(\mathbf{y}^n)$ such that $aC \leq X_n \leq bC$ and use the definitions of $\mathbb{H}(p(\mathbf{y}))$ and $\hat{\mathbb{H}}(p(\mathbf{y}))$, obtaining

$$\mathbb{P} \left(\left| \hat{\mathbb{H}}(p(\mathbf{y})) - \mathbb{H}(p(\mathbf{y})) \right| \geq \epsilon \right) \leq 2 \exp \left(-\frac{2N\epsilon^2}{(b-a)^2 C^2} \right). \quad (18)$$

We set the r.h.s. $\leq \delta$ and solve for N to arrive at the desired bound.

Result. To ensure with probability at least $1 - \delta$ that $\left| \hat{\mathbb{H}}(p(\mathbf{y})) - \mathbb{H}(p(\mathbf{y})) \right| \leq \epsilon$, it suffices that

$$N \geq \frac{(b-a)^2 C^2}{2\epsilon^2} \log \left(\frac{2}{\delta} \right) \quad (19)$$

Next, we discuss the practical implications of this bound and contrast it to Eq. (16).

B.3 COMPARING THEORETICAL APPROXIMATION QUALITY

The bounds in Eq. (16) and Eq. (19) quantify the number of samples needed to achieve a given approximation error ϵ with high probability $1 - \delta$. By analyzing how the bounds scale with properties of p and the sampling distribution q , we can compare the relative difficulty of estimating each quantity.

The bound for $\hat{\mathbb{M}}(p(\mathbf{y}))$ in Eq. (16) depends only on the concentration of the output sequence distribution along a few probable sequences within the ϵ -region. This is amenable to decoding practices in LLMs such as top-k, top-p, low-temperature, or even beam search, that focus on obtaining very likely output sequences.

The bound for $\hat{\mathbb{H}}(p(\mathbf{y}))$ in Eq. (19), in contrast, reflects the variance of the entire distribution and depends on the squared range of $\log p(\mathbf{y})$ scaled by the worst-case importance weight C^2 , both of which can be very high in practice. As it aggregates over the entire support of $p(\mathbf{y})$, estimation is sensitive to rare but non-negligible sequences that are unlikely to be sampled under practical sampling strategies.

Regarding importance sampling, approximating $\mathbb{M}(p(\mathbf{y}))$ is also desirable. Here, it suffices that the proposal distribution q concentrates around the most likely sequences according to p and there are no importance weights in the calculation of $\hat{\mathbb{M}}(p(\mathbf{y}))$. Importance sampling only increases the likelihood of sampling close to the maximum, not the calculation of the estimator itself. In contrast, when approximating $\mathbb{H}(p(\mathbf{y}))$, it is necessary to re-weight the samples using importance weights $c(\mathbf{y})$. Under a less optimal q , this can lead to a strong increase in variance.

Overall, from a sample complexity point of view, approximating $\mathbb{M}(p(\mathbf{y}))$ is not only more efficient than approximating $\mathbb{H}(p(\mathbf{y}))$, but also well aligned with the way LLMs are used in practice.

C DETAILS ON SIMULATION STUDY

In this section, we provide detailed insights into the numerical results presented in Sec. 3, especially Fig. 1.

To recall, we empirically investigate the performance of estimators for the predictive entropy $H(p(\mathbf{y}))$ (Eq. (4)) and the negative maximum sequence log-likelihood (min-entropy) $M(p(\mathbf{y}))$ (Eq. (9)) in a controlled setting. Therefore, we consider a synthetic experiment with the following setup. We are given a space of possible outcomes \mathcal{V} with $|\mathcal{V}| = \{20, 100\}$. The task is to predict a sequence $\mathbf{y} = (y_1, \dots, y_T) \in \mathcal{V}_T$ where $y \in \mathcal{V}$ and T is 2, 3, or 4. Predictive distributions $p(\mathbf{y})$ are not represented by a neural network, but are randomly sampled according to a Dirichlet distribution $\text{Dir}(\{\alpha_1, \dots, \alpha_{|\mathcal{V}|}\})$. The alpha parameters of the Dirichlet distribution are specified to yield typical predictive distributions as encountered in LLMs that follow a Zipf distribution. For $|\mathcal{V}| = 20$ we have $\alpha_{1,2} = 10$ and $\alpha_{3-20} = 0.2$. For $|\mathcal{V}| = 100$ we have $\alpha_{1,2} = 10$, $\alpha_{3-6} = 1$ and $\alpha_{7-100} = 0.2$. Note that the order of alpha values is randomly shuffled before drawing each predictive distribution. Representative predictive distributions sampled from this Dirichlet distribution are shown in Fig. 3a and Fig. 3b.

The experiments investigate the quality of the estimators depending on the number of samples $\{\mathbf{y}_n\}_{n=1}^N$. This is feasible because the ground truth values for both entropy and maximum sequence log-likelihood can be calculated for this small synthetic example through exhaustive enumeration. For the experiments we present here in the appendix, we average over 2,000 runs, meaning that new sets of samples $\{\mathbf{y}_n\}_{n=1}^N$ are drawn to calculate the respective estimators. As beam search is deterministic, it does not vary in this experimental setting, compared to Fig. 1b in the main paper, where we investigated the quality of estimators for different $p(\mathbf{y})$.

The results for estimating the entropy are shown in Fig. 4. We observe that the variance of estimators increases for larger vocabulary sizes $|\mathcal{V}|$ and sequence lengths T . Furthermore, lower temperatures increase the variance of the estimator as expected. Note that we introduced clipping for importance weights at $1e \pm 6$, which did not introduce noticeable bias, but made results numerically stable.

The results for estimating the maximum sequence likelihood are shown in Fig. 5. We observe that low-temperature multinomial sampling and beam search find the maximum sequence log-likelihood with a low number of samples with high probability. Greedy decoding (beam width of 1) finds the maximum for all experimental settings except one, where it takes a beam width of 2 to find it with high probability.

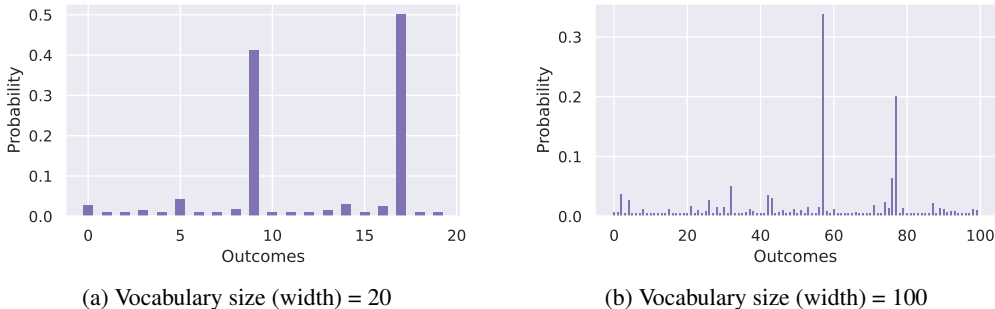


Figure 3: Exemplary predictive distributions $p(y_t | \mathbf{y}_{<t})$ for different vocabulary sizes (widths).

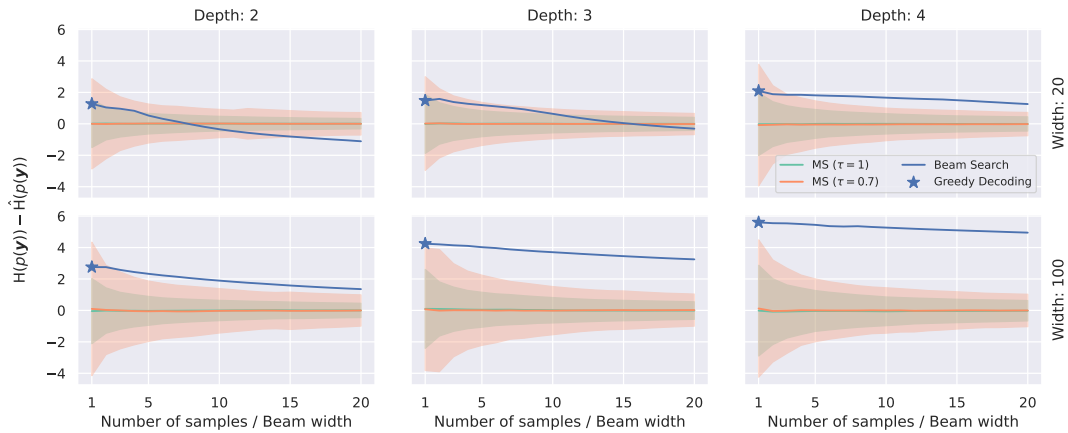


Figure 4: **Estimator of Predictive Entropy.** Results for different vocabulary sizes (width) and sequence lengths (depth). We estimate the entropy $H(p(\mathbf{y}))$ using N Monte-Carlo samples (cf. Eq. (4)). Lines denote the average over runs, while shades denote one standard deviation. We compare MS for two commonly used τ . The experiments show that the decreased temperature leads to lower variance but introduces bias.

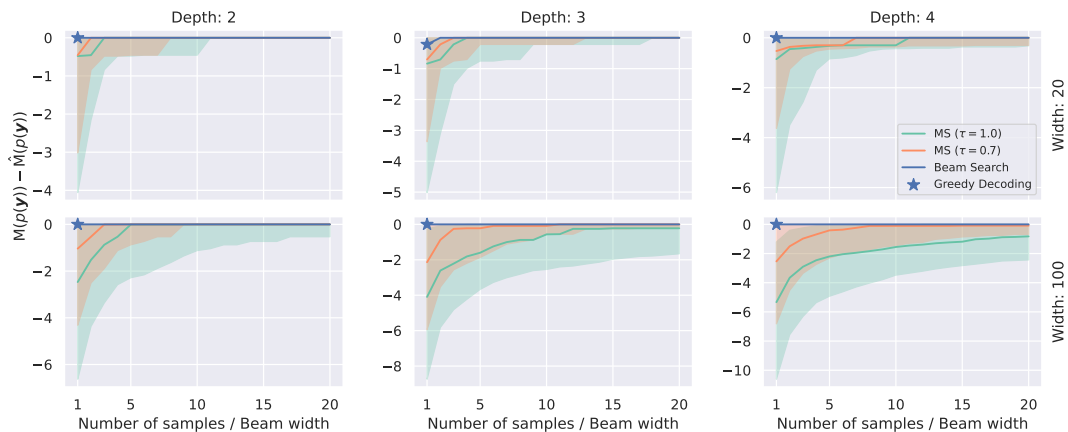


Figure 5: **Estimator of maximum sequence log-likelihood.** Results for different vocabulary sizes (width) and sequence lengths (depth). We estimate $M(p(\mathbf{y}))$ using the maximum over N sampled obtained by beam search ($N = 1$ is greedy decoding) or MS with different τ . Lines denote the median, and shades signify the possible values between the 5 and 95 percent quantiles. Beam search is deterministic for any given distribution $p(\mathbf{y})$. Even with a very low number of samples, low-temperature MS and beam search are able to find the maximum with high probability.

D DETAILS ON EXPERIMENTAL RESULTS

In this section, we provide detailed insights to complement the main results presented in Sec. 5.1. The code and data are available at <https://github.com/ml-jku/G-NLL>.

Hyperparameters. To compute G-NLL, we use greedy decoding to generate the reference answer, which is equivalent to beam search with a single beam and multinomial sampling with a sampling temperature of zero. Notably, G-NLL is deterministic and hyperparameter-free. To compute the logarithmic score based measures (PE , $LN-PE$, SE , $LN-SE$, $D-SE$), 10 output sequences are generated via multinomial sampling. For each dataset, we performed a hyperparameter search on held-out instances to determine the best-performing temperature $\tau \in \{0.5, 1.0, 1.5\}$ for sampling output sequences used for the logarithmic score based measures.

Evaluation Metrics. The correctness of the reference answer is assessed by checking if the F1 score of the commonly used SQuAD metric between the reference answer and the given answer exceeds 0.5 ($F1$). Although there are some limitations to using such a simple metric, it has relatively small errors in standard datasets and, therefore, remains widely used in practice. Additionally, we use the model Llama-3.1-70B as LLM-as-a-judge to assess if the given answer is correct (LLM).

AUROC is used as the primary performance metric throughout this paper, consistent with standard evaluation practices in this field. We report the results for the individual datasets in Tab. 2, which have been averaged over in Tab. 1. In addition to AUROC, we also consider the average rejection accuracy, i.e., the accuracy of model predictions when allowing the rejection of a certain budget of predictions based on the uncertainty estimate. Results are presented in Tab. 3, where predictions are only evaluated for the 80% most certain predictions, and we again use greedy decoding for our measure based on the zero-one score. This further suggests that our measure remains highly competitive across various settings.

Evaluation Tasks. We evaluate our method on both short phrase answers and full-sentence answers (referred to as “short” and “long” respectively), assessing the generalization of uncertainty estimation across model classes, model sizes, training stages, and evaluation criteria for output correctness, using three different question-answering datasets. We selected these datasets to best align with prior work to provide a fair and meaningful comparison (Duan et al., 2024; Kuhn et al., 2023; Farquhar et al., 2024; Bakman et al., 2024; Nikitin et al., 2024; Aichberger et al., 2025; Kossen et al., 2024). Notably, Farquhar et al. (2024) additionally investigates paragraph-length summarization of biographies by decomposing the task into individual question-answering problems. This suggests that performance on question-answering tasks is expected to be indicative of performance on summarization tasks as well. As such, our experimental setup reflects a wide spectrum of practical tasks.

Length-Normalization. Intuitively, length normalization might appear to be a reasonable choice also for the single-sequence measure. However, we find that G-NLL without length normalization yields the best performance, even on tasks with high variability in sequence length, as reported in Tab. 4. This can be attributed to the fact that length normalization tends to dilute the influence of low-probability tokens. Since most tokens typically have relatively high likelihoods, summing log-probabilities (rather than averaging them) places greater emphasis on the more uncommon, low-likelihood tokens that are more informative for uncertainty estimation. Although not specific to single-sequence measures, investigating the trade-off introduced by length normalization to handle these characteristics represents a promising direction for future work.

Beam-Search for MSP Approximation. We investigate how much the alternative measure derived from the zero-one score benefits from better approximating the most likely output sequences (MSP), as searching for more likely output sequences through beam search theoretically further improves the approximation of MSP. The results summarized in Tab. 5 show that the performance does not significantly improve compared to greedy decoding, which is consistent with the ablation study presented in Sec. 5.2. This further supports the claim that G-NLL is a strong measure of uncertainty, despite its algorithmic simplicity and computational efficiency.

Table 2: **AUROC Evaluation.** Average AUROC (\uparrow) using uncertainty estimates of different measures as a score to distinguish between correct and incorrect answers of each dataset. The reference answer is generated using greedy decoding, either as a whole sentence (*long*) or a short phrase (*short*).

\mathcal{D}	Uncertainty measure generating scoring rule				Logarithmic				Zero-One			
	Language Model	Gen.	Metric	PE	LN-PE	SE	LN-SE	D-SE	G-NLL			
TriviaQA	Transformer	8B	short	F1	.758	.778	.788	<u>.798</u>	.787	.810		
			PT	short	LLM	.675	.694	.703	<u>.704</u>	.682	.722	
			long	LLM	.592	.604	.640	.631	<u>.650</u>	.704		
			IT	short	F1	.735	.768	.790	<u>.800</u>	.777	.809	
			PT	short	LLM	.660	.684	.708	<u>.710</u>	.680	.716	
			long	LLM	.603	.627	.678	.672	.670	.670		
		70B	short	F1	.707	.730	.741	<u>.743</u>	.702	.744		
			PT	short	LLM	.650	.660	<u>.696</u>	.695	.656	.698	
			long	LLM	.538	.533	<u>.625</u>	.574	.563	.692		
			IT	short	F1	.698	.714	.722	.726	.688	.722	
			PT	short	LLM	.663	.675	<u>.685</u>	.679	.633	.701	
			long	LLM	.530	.553	.564	.571	.564	.543		
	State-Space	7B	short	F1	.786	.793	.812	<u>.818</u>	.810	.832		
			PT	short	LLM	.687	.697	.712	<u>.714</u>	.695	.724	
			long	LLM	.597	.653	.675	<u>.680</u>	<u>.689</u>	.705		
			IT	short	F1	.780	.799	.810	<u>.819</u>	.811	.827	
			PT	short	LLM	.696	.701	.714	<u>.717</u>	.703	.730	
			long	LLM	.645	.654	.688	.698	.692	.694		
	SVAMP	Transformer	8B	short	F1	.847	.867	.865	<u>.870</u>	.868	.885	
				PT	short	LLM	.779	.788	.753	<u>.772</u>	.791	.772
				long	LLM	.575	.563	.519	.534	<u>.601</u>	.669	
				IT	short	F1	.879	.903	<u>.914</u>	.912	<u>.887</u>	.931
				PT	short	LLM	.706	.725	<u>.736</u>	.731	.701	.753
				long	LLM	.556	.524	.590	.608	<u>.631</u>	.662	
70B			short	F1	.892	.906	.925	<u>.929</u>	.923	.936		
			PT	short	LLM	.794	.817	.814	.815	.819	.799	
			long	LLM	.578	.554	.553	<u>.579</u>	.571	.665		
			IT	short	F1	.830	.895	.915	.922	.915	.909	
			PT	short	LLM	.703	.744	.734	.748	.762	.713	
			long	LLM	.601	.577	.613	.649	.663	.597		
State-Space		7B	short	F1	.882	<u>.893</u>	.874	.883	.889	.914		
			PT	short	LLM	.752	<u>.757</u>	.730	.738	<u>.757</u>	.776	
			long	LLM	.536	.585	.534	.602	.612	.579		
			IT	short	F1	.843	.891	.854	.876	<u>.892</u>	.905	
			PT	short	LLM	.706	.730	.704	.709	<u>.737</u>	.744	
			long	LLM	.577	.586	.578	.616	.639	.613		
NQ		Transformer	8B	short	F1	.725	.739	.673	.710	<u>.758</u>	.776	
				PT	short	LLM	.639	.661	.615	.641	.683	.683
				long	LLM	.517	.498	.478	.495	<u>.550</u>	.573	
				IT	short	F1	.702	.732	.711	.731	<u>.756</u>	.774
				PT	short	LLM	.662	.682	.669	.685	.700	.697
				long	LLM	.494	.491	.530	.524	.527	.514	
	70B		short	F1	.727	.733	.711	.737	<u>.748</u>	.779		
			PT	short	LLM	.634	.649	.642	.657	<u>.671</u>	.672	
			long	LLM	.538	.514	.494	.553	<u>.580</u>	.589		
			IT	short	F1	.718	.734	.734	.748	.746	.743	
			PT	short	LLM	.676	.674	.689	.698	.702	.681	
			long	LLM	.535	.540	.526	.566	.574	.545		
	State-Space	7B	short	F1	.766	.758	.741	.765	.785	.782		
			PT	short	LLM	.675	.680	.661	.681	.697	.683	
			long	LLM	.567	.553	.512	.551	.572	.554		
			IT	short	F1	.755	.751	.727	.754	.783	.781	
			PT	short	LLM	.669	.672	.648	.671	.692	.683	
			long	LLM	.541	.521	.526	.541	.554	.537		
	Average				.677	.689	.690	.703	.707	.721		

Table 3: **Rejection Accuracy Evaluation.** Average Rejection Accuracy at 80% (\uparrow) across all datasets, using uncertainty estimates of different measures as a score to distinguish between correct and incorrect answers. The reference answer is generated using greedy decoding, either as a whole sentence (*long*) or a short phrase (*short*).

<i>Uncertainty measure generating scoring rule</i>			<i>Logarithmic</i>					<i>Zero-One</i>	
Language Model	Gen.	Metric	PE	LN-PE	SE	LN-SE	D-SE	G-NLL	
Transformer	8b	short	F1	.661	<u>.672</u>	.651	.643	.655	.681
			LLM	.774	.782	.767	.766	.765	.778
			LLM-Instruct	.704	<u>.721</u>	.693	.688	.702	.723
		long	LLM	.596	<u>.590</u>	<u>.598</u>	.592	.590	.619
			LLM-Instruct	.667	<u>.684</u>	.632	.643	.644	.686
			F1	.668	.684	.680	.673	<u>.687</u>	.702
	IT	short	LLM	.775	<u>.781</u>	.779	.775	.778	.788
			LLM-Instruct	.723	.742	.732	.726	<u>.743</u>	.751
			LLM	.628	.630	.651	.644	.653	.652
		long	LLM-Instruct	.713	.724	.705	.713	<u>.727</u>	.734
			F1	.818	.827	.822	.827	<u>.829</u>	.836
			LLM	.844	<u>.852</u>	.846	.847	.851	.855
70b	short	LLM-Instruct	.867	.875	.876	.881	.885	.881	
		LLM	.704	.699	<u>.719</u>	.707	.705	.724	
		LLM-Instruct	.789	<u>.795</u>	.776	.781	.788	.812	
	long	F1	.795	.813	.814	.809	<u>.819</u>	.823	
		LLM	.836	.842	.842	.837	<u>.844</u>	.845	
		LLM-Instruct	.850	.867	.866	.865	.874	.870	
7b	short	LLM	.706	.706	.712	.715	.721	.715	
		LLM-Instruct	.855	.850	.827	.842	.861	.851	
		F1	<u>.598</u>	.596	.585	.579	.583	.612	
	long	LLM	.729	<u>.737</u>	.723	.721	.733	.742	
		LLM-Instruct	.638	<u>.640</u>	.626	.621	.632	.651	
		LLM	.613	.627	.612	.624	.620	.623	
State-Space	short	LLM-Instruct	.606	.611	.601	.611	<u>.618</u>	.633	
		F1	.592	<u>.603</u>	.588	.581	.589	.615	
		LLM	.737	<u>.742</u>	.730	.726	.740	.744	
	long	LLM-Instruct	.632	<u>.646</u>	.625	.619	.637	.653	
		LLM	.611	.617	.618	.612	.625	.625	
		LLM-Instruct	.643	.652	.628	.628	<u>.654</u>	.658	
Average			.712	.720	.711	.710	.718	.729	

Table 4: **Length Normalization Evaluation.** Average AUROC (\uparrow) across TriviaQA, SVAMP, and NQ datasets, utilizing G-NLL and its length-normalized version LN-G-NLL to assign an uncertainty estimate, which is used as a score to distinguish between correct and incorrect answers. The reference answer is generated using greedy decoding, either as a whole sentence (*long*) or a short phrase (*short*).

			<i>Uncertainty measure generating scoring rule</i>		<i>Zero-One</i>	
	Language Model	Generation	Metric	LN-G-NLL	G-NLL	
Transformer	8B	PT	short	F1	.811	.824
			short	LLM	.717	.726
			long	LLM	.532	.649
		IT	short	F1	.826	.838
			short	LLM	.716	.722
			long	LLM	.542	.615
	70B	PT	short	F1	.811	.820
			short	LLM	.718	.723
			long	LLM	.529	.649
		IT	short	F1	.788	.792
			short	LLM	.695	.699
			long	LLM	.539	.562
State-Space	7B	PT	short	F1	.828	.843
			short	LLM	.720	.728
		long	LLM	.576	.612	
		IT	short	F1	.823	.838
	short		LLM	.713	.719	
	long	LLM	.562	.615		

Table 5: **Beam Search Evaluation.** Average AUROC (\uparrow) across all datasets, using uncertainty estimates of different measures as a score to distinguish between correct and incorrect answers. The reference answer is generated using *beam search with 5 beams*, again either as a whole sentence (*long*) or a short phrase (*short*).

			<i>Uncertainty measure generating scoring rule</i>			<i>Logarithmic</i>			<i>Zero-One</i>	
	Language Model	Gen.	Metric	PE	LN-PE	SE	LN-SE	D-SE	G-NLL	
Transformer	8B	PT	short	F1	.775	.791	.765	.787	<u>.799</u>	.822
			short	LLM	.700	.712	.686	.704	<u>.713</u>	.726
			long	LLM	.556	.540	.493	.520	<u>.578</u>	.591
		IT	short	F1	.778	.808	.805	<u>.819</u>	.811	.845
			short	LLM	.682	.704	.706	<u>.713</u>	.698	.729
			long	LLM	.535	.520	.584	.585	.586	.559
	70B	PT	short	F1	.788	.799	.796	<u>.812</u>	.798	.833
			short	LLM	.700	.717	.719	.727	.718	.725
			long	LLM	.540	.552	.489	.531	<u>.552</u>	.608
		IT	short	F1	.756	.786	.796	.806	.788	.800
			short	LLM	.680	.697	.701	.707	.695	.707
			long	LLM	.534	.533	.544	.569	.574	.534
State-Space	7b	PT	short	F1	.814	.818	.806	.823	<u>.825</u>	.846
			short	LLM	.703	.709	.699	.711	<u>.712</u>	.719
		long	LLM	.570	.595	.550	.609	.602	.563	
		IT	short	F1	.799	.815	.794	.817	<u>.828</u>	.845
	short		LLM	.699	.713	.694	.709	<u>.720</u>	.730	
	long	LLM	.574	.575	.582	.621	.607	.577		
	Average				.677	.688	.678	.698	.700	.709