

# RETHINKING UNCERTAINTY ESTIMATION IN NATURAL LANGUAGE GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) are increasingly employed in real-world applications, driving a need to determine when their generated text can be trusted or should be questioned. To assess the trustworthiness of the generated text, reliable uncertainty estimation is essential. Current LLMs generate text through a stochastic process that can lead to different output sequences for the same prompt. Consequently, leading uncertainty measures require generating multiple output sequences to estimate the LLM’s uncertainty. However, generating additional output sequences is computationally expensive, making these uncertainty estimates impractical at scale. In this work, we challenge the theoretical foundations of the leading measures and derive an alternative measure that eliminates the need for generating multiple output sequences. Our new measure is based solely on the negative log-likelihood of the most likely output sequence. This vastly simplifies uncertainty estimation while maintaining theoretical rigor. Empirical results demonstrate that our new measure achieves state-of-the-art performance across various models and tasks. Our work lays the foundation for reliable and efficient uncertainty estimation in LLMs, challenging the necessity of the more complicated methods currently leading the field.

## 1 INTRODUCTION

Language models are increasingly adopted in a wide range of real-world applications. Despite the advancements in language models, determining whether a generated text can be trusted remains a significant challenge. To address this challenge, it is crucial to reliably assess the level of certainty a language model has regarding its generated text. While uncertainty estimates do not guarantee factuality for generated text based on consistent but erroneous training data, they are a reliable indicator of errors at present (Kuhn et al., 2023; Aichberger et al., 2024; Farquhar et al., 2024).

Assessing predictive uncertainty in language models is inherently difficult due to their autoregressive nature. For a given input sequence, language models predict the next token probabilities, based on which a specific token is selected and appended to the sequence. This stochastic process is repeated for each new token. Selecting different tokens at specific steps during the generation leads to varying output sequences for the same input sequence with the same language model. Consequently, the space of possible output sequences is vast and computationally intractable to fully explore (Sutskever et al., 2014; Vaswani et al., 2017; Radford et al., 2018).

Current uncertainty estimation methods rely on assessing the probability distribution over all possible output sequences. However, the generation of each additional token is computationally expensive, and practical methods can only sample a small fraction of possible output sequences (Malinin & Gales, 2021; Kadavath et al., 2022). Moreover, even after having generated multiple likely output sequences, the question remains whether these indicate high uncertainty. A language model that likely generates different output sequences is not necessarily uncertain about the underlying meaning if the output sequences are semantically equivalent. Leading uncertainty measures address this fact by considering the semantics of the output sequences, utilizing separate language inference models (Kuhn et al., 2023; Farquhar et al., 2024). While these measures improve the performance of the uncertainty estimates, they also further add complexity and computational overhead. These factors make current uncertainty estimation methods impractical at scale, hindering their broad adoption in

054 real-world applications. There is a need for efficient uncertainty estimation methods that give clear  
 055 insights into the reliability of language models without incurring substantial computational costs.

056  
 057 In this work, we assess whether we can theoretically motivate an uncertainty measure that does not  
 058 rely on the probability distribution over all possible output sequences. Building on insights from  
 059 the principled framework of proper scoring rules (Kotelevskii & Panov, 2024; Hofman et al., 2024),  
 060 we adopt the zero-one score as an alternative to the currently used logarithmic score for uncertainty  
 061 measures in NLG. This leads to a theoretically motivated measure that does not require generating  
 062 multiple output sequences but solely relies upon a single output sequence. Our proposed measure  
 063 is straightforward: it simply is the negative log-likelihood of the most likely output sequence. By  
 064 eliminating the need to generate and semantically cluster multiple output sequences, our measure  
 065 significantly reduces computational costs and complexity.

066  
 067 Experimental results demonstrate that our new measure matches or even exceeds the performance  
 068 of current state-of-the-art uncertainty estimation methods across various model classes, model sizes,  
 069 model stages, tasks, datasets, and evaluation metrics. In summary, our new measure not only pre-  
 070 serves theoretical rigor but also provides a more scalable solution for uncertainty estimation in lan-  
 071 guage models, making it highly practical for real-world applications.

072 Our main contributions are:

- 073 • We introduce the negative log-likelihood of the most likely output sequence as an efficient and  
 074 practical measure of uncertainty in NLG.
- 075 • We provide a rigorous theoretical foundation for our measure, building upon established principles  
 076 in uncertainty theory and proper scoring rules.
- 077 • We conduct extensive experiments demonstrating that our measure achieves strong performance,  
 078 matching or surpassing state-of-the-art methods while significantly reducing computational costs.

## 079 2 PREDICTIVE UNCERTAINTY IN NLG

080  
 081 **Preliminaries.** We assume a fixed training dataset  $\mathcal{D} = \{s_i\}_{i=1}^N$  consisting of ordered tokens  
 082  $s_t \in \mathcal{V}$ , with  $\mathcal{V}$  being a given vocabulary. Each token at step  $t$  is assumed to be sampled according  
 083 to the predictive distribution  $p(s_t | s_{<t}, w^*)$ , conditioned on the sequence of preceding tokens  
 084  $s_{<t}$  and the true (but unknown) language model parameters  $w^*$ . We assume that the given model  
 085 class can theoretically represent the true predictive distribution, a common and usually necessary  
 086 assumption (Hüllermeier & Waegeman, 2021). How likely language model parameters  $\tilde{w}$  match  
 087  $w^*$  is determined by the posterior distribution  $p(\tilde{w} | \mathcal{D}) = p(\mathcal{D} | \tilde{w})p(\tilde{w})/p(\mathcal{D})$ .

088 In language model inference, the input to a given language model parameterized by  $w$  is a sequence  
 089  $\mathbf{x} = (x_1, \dots, x_M)$  and the output is a sequence  $\mathbf{y} = (y_1, \dots, y_T) \in \mathcal{Y}_T$ , with  $x, y \in \mathcal{V}$  and  $\mathcal{Y}_T$   
 090 being the set of all possible output sequences with a sequence length smaller equal to  $T$ . The  
 091 likelihood of a token  $y_t \in \mathbf{y}$  being generated by the language model is conditioned on both the input  
 092 sequence and all previously generated tokens, denoted as  $p(y_t | \mathbf{x}, \mathbf{y}_{<t}, w)$ . The likelihood of output  
 093 sequences  $\mathbf{y} \in \mathcal{Y}_T$  being generated by the language model is then the product of the individual token  
 094 probabilities, denoted as  $p(\mathbf{y} | \mathbf{x}, w) = \prod_{t=1}^T p(y_t | \mathbf{x}, \mathbf{y}_{<t}, w)$  (Sutskever et al., 2014), while the  
 095 heuristic length-normalized variant is  $\bar{p}(\mathbf{y} | \mathbf{x}, w) = \exp\left(\frac{1}{T} \sum_{t=1}^T \log p(y_t | \mathbf{x}, \mathbf{y}_{<t}, w)\right)$  (Malinin  
 096 & Gales, 2021).

097 Computing the likelihood of a specific output sequence  $\mathbf{y}$  being generated by the language model  
 098 parameterized by  $w$  – or in other words, being sampled from the probability distribution over pos-  
 099 sible output sequences  $\mathbf{y} \sim p(\mathbf{y} | \mathbf{x}, w)$  – is straightforward. The language model directly provides  
 100 the individual token likelihoods. However, determining the full probability distribution over pos-  
 101 sible output sequences is considerably more challenging, since  $\mathcal{Y}_T$  scales exponentially with the  
 102 sequence length  $T$ . The computational complexity of evaluating all possible sequences grows as  
 103  $\mathcal{O}(|\mathcal{V}|^T)$ . Since modern language models even exceed a vocabulary size  $|\mathcal{V}|$  of one hundred thou-  
 104 sand tokens, this distribution becomes intractable to compute, even for relatively short maximal  
 105 sequence lengths  $T$  (Dubey et al., 2024).

**Uncertainty Measures and Proper Scoring Rules.** We now derive measures to estimate uncertainty in NLG. Throughout this work, the focus is on estimating the predictive uncertainty of a single, given “off-the-shelf” model. We assume to that a given language model parameterized by  $w$  is the *predicting model* used to sample output sequences  $\mathbf{y} \sim p(\mathbf{y} \mid \mathbf{x}, w)$ . Furthermore, we assume that any language model parameterized by  $\tilde{w}$  is an *approximation of the true predictive distribution* according to its posterior probability  $p(\tilde{w} \mid \mathcal{D})$ . Together, these two assumptions give rise to specific uncertainty measures (Schweighofer et al., 2023; 2024), as elaborated on in more detail below. Aichberger et al. (2024) shows that established uncertainty measures in NLG, such as Predictive Entropy (PE) (Malinin & Gales, 2021) and Semantic Entropy (SE) (Kuhn et al., 2023; Farquhar et al., 2024), naturally emerge under this assumption. In general, the information-theoretic entropy has become the standard measure to assess predictive uncertainty. However, recent studies by Lahlou et al. (2023); Gruber & Buettner (2023); Kotelevskii & Panov (2024) and Hofman et al. (2024) have shown that these information-theoretic measures are not the only viable options. A broader class of *proper scoring rules* provides a principled framework for predictive uncertainty measures. In the following, we leverage this framework to derive our alternative measure that relies solely on a single output sequence. We begin by discussing the concept of proper scoring rules.

In general, 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 the actual observations (Gneiting & Raftery, 2007). For uncertainty estimation in NLG, the general notion of proper scoring rules assigns a numerical score to how well the predicted distribution of output sequences  $p(\mathbf{y} \mid \mathbf{x}, \cdot)$  aligns with the observed output sequence  $\mathbf{y}'$ . In particular, a proper scoring rule is an extended real-valued function  $S : \mathcal{P} \times \mathcal{Y} \rightarrow [-\infty, \infty]$ , such that  $S(p, \cdot)$  is  $\mathcal{P}$ -quasi-integrable over a convex class of probability measures  $\mathcal{P}$ . The expected score over possible output sequences  $\mathbf{y}'$  is given by

$$\mathbb{E}_{\mathbf{y}' \sim p(\mathbf{y}' \mid \mathbf{x}, \cdot)} [S(p(\mathbf{y} \mid \mathbf{x}, \cdot), \mathbf{y}')] \quad (1)$$

Given this general formulation, we now incorporate the assumptions outlined above to establish the connection to uncertainty measures (Schweighofer et al., 2024). First, under the assumption about the *predicting model*, the distribution giving rise to the observed output sequences  $p(\mathbf{y}' \mid \mathbf{x}, \cdot)$  corresponds to the predictive distribution of the given language model, denoted as  $p(\mathbf{y}' \mid \mathbf{x}, w)$ . Second, under the assumption about the *approximation of the true predictive distribution*, a sampled output sequence  $\mathbf{y}'$  has to be compared to all possible language models parameterized by  $\tilde{w}$ , according to their posterior distribution  $p(\tilde{w} \mid \mathcal{D})$ . This captures how much the sampled output sequence aligns with all possible predictive distributions  $p(\mathbf{y} \mid \mathbf{x}, \tilde{w})$ . Therefore, we take a posterior expectation over Eq. (1), which can be additively decomposed into an entropy term and a divergence term (Gneiting & Raftery, 2007; Kull & Flach, 2015):

$$\begin{aligned} & \mathbb{E}_{\tilde{w} \sim p(\tilde{w} \mid \mathcal{D})} \left[ \underbrace{\mathbb{E}_{\mathbf{y}' \sim p(\mathbf{y}' \mid \mathbf{x}, w)} [S(p(\mathbf{y} \mid \mathbf{x}, \tilde{w}), \mathbf{y}')] }_{\text{expected score}} \right] \quad (2) \\ & = \underbrace{\mathbb{E}_{\mathbf{y}' \sim p(\mathbf{y}' \mid \mathbf{x}, w)} [S(p(\mathbf{y} \mid \mathbf{x}, w), \mathbf{y}')] }_{\text{entropy term}} \\ & \quad + \underbrace{\mathbb{E}_{\tilde{w} \sim p(\tilde{w} \mid \mathcal{D})} \left[ \mathbb{E}_{\mathbf{y}' \sim p(\mathbf{y}' \mid \mathbf{x}, w)} [S(p(\mathbf{y} \mid \mathbf{x}, \tilde{w}), \mathbf{y}') - S(p(\mathbf{y} \mid \mathbf{x}, w), \mathbf{y}')] \right]}_{\text{divergence term}} . \end{aligned}$$

**Aleatoric and Epistemic Uncertainty.** In terms of predictive uncertainty, this general framework can be interpreted as follows. The expected score over possible output sequences and language model parameters captures the *total uncertainty of the given language model*. The entropy term reflects *aleatoric* uncertainty, which is the uncertainty inherent in the data generation process, arising from the inherent variability and randomness in natural language (Gal, 2016; Kendall & Gal, 2017). 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 or model capacity (Houlsby et al., 2011; Gal, 2016; Malinin, 2019; Hüllermeier & Waegeman, 2021).

The concrete *total*, *aleatoric*, and *epistemic* uncertainty measures depends on the choice of proper scoring rule. For instance, the logarithmic score is the most common proper scoring rule that gives rise to the well-known information-theoretic uncertainty measures in both classification tasks (Houlsby et al., 2011; Gal, 2016) and NLG (Malinin & Gales, 2021; Kuhn et al., 2023).

In the following, we first revisit these uncertainty measures that are based on the logarithmic score and analyze their effectiveness in estimating aleatoric and epistemic uncertainty. Thereafter, we propose uncertainty measures that are based on another proper scoring rule, the zero-one score. This score has not yet been considered for uncertainty estimation in NLG. We show that utilizing uncertainty measures based on the zero-one score offers certain advantages.

## 2.1 ESTABLISHED UNCERTAINTY MEASURES IN NLG BASED ON LOGARITHMIC SCORE

The logarithmic score is usually assumed implicitly to derive uncertainty measures, due to the foundation of resulting measures in information theory (Lahlou et al., 2023; Gruber & Buettner, 2023; Hofman et al., 2024; Kotelevskii & Panov, 2024). In the context of NLG, it considers the negative log-likelihood of a generated output sequence  $\mathbf{y}'$ :

$$S_{\log}(p(\mathbf{y} | \mathbf{x}, \cdot), \mathbf{y}') = -\log p(\mathbf{y} = \mathbf{y}' | \mathbf{x}, \cdot). \quad (3)$$

Using the logarithmic score in Eq. (2) results in the cross-entropy  $\text{CE}(\cdot; \cdot)$  between the output sequence distribution of the given language model and that of every possible language model according to their posterior  $p(\tilde{\mathbf{w}} | \mathcal{D})$  (Schweighofer et al., 2023; Aichberger et al., 2024):

$$\begin{aligned} & \underbrace{E_{\tilde{\mathbf{w}} \sim p(\tilde{\mathbf{w}} | \mathcal{D})} [\text{CE}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}); p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{total}} \\ &= \underbrace{H(p(\mathbf{y} | \mathbf{x}, \mathbf{w}))}_{\text{aleatoric}} + \underbrace{E_{\tilde{\mathbf{w}} \sim p(\tilde{\mathbf{w}} | \mathcal{D})} [\text{KL}(p(\mathbf{y} | \mathbf{x}, \mathbf{w}) \| p(\mathbf{y} | \mathbf{x}, \tilde{\mathbf{w}}))]}_{\text{epistemic}}. \end{aligned} \quad (4)$$

The epistemic uncertainty is a posterior expectation of the Kullback-Leibler divergence  $\text{KL}(\cdot \| \cdot)$  between the output sequence distribution of the given model and that of all possible models. This requires considering every possible model parametrization. Since modern language models have billions of parameters (Radford et al., 2018; Zhang et al., 2022; Touvron et al., 2023; Zuo et al., 2024; Dubey et al., 2024), the epistemic uncertainty is particularly challenging to estimate.

Current work usually solely considers the aleatoric uncertainty, which is the Shannon entropy  $H(\cdot)$  of the output sequence distribution of the given language model (Malinin & Gales, 2021; Kuhn et al., 2023; Aichberger et al., 2024). Computing the output sequence distribution still requires considering the whole set of possible output sequences  $\mathcal{Y}_T$ . Thus, the primary objective of uncertainty estimation based on the logarithmic score is to closely approximate this output sequence distribution.

**Predictive Entropy.** The aleatoric uncertainty under a given language model is the entropy of the output sequence distribution, commonly referred to as Predictive Entropy (PE). Intuitively, high PE implies that the language model is likely to generate different output sequences from the same input sequence, indicating high uncertainty of the language model. PE usually is estimated via Monte Carlo (MC) sampling (Malinin & Gales, 2021):

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

**Semantic Entropy.** Semantic Entropy (SE) builds on the fact that output sequences may be different on a token level but equivalent on a semantics level. In such cases, the PE can be misleading, as it reflects high uncertainty even when different output sequences have the same semantic meaning. PE also captures the uncertainty of the language model in expressing the semantically same statement, which is often not the focus of uncertainty estimation in NLG. Thus, instead of the entropy of the output sequence distribution, the entropy of the semantic cluster distribution is considered, denoted as  $p(c | \mathbf{x}, \mathbf{w}) = \sum_{\mathbf{y}} p(c | \mathbf{x}, \mathbf{y}, \mathbf{w}) p(\mathbf{y} | \mathbf{x}, \mathbf{w})$ . The probability of an output sequence belonging to a semantic cluster is usually approximated with a separate natural language inference model. High SE implies that the language model is likely to generate output sequences that have high semantic diversity, indicating high semantic uncertainty (Kuhn et al., 2023; Farquhar et al., 2024).

$$\begin{aligned} H(p(c | \mathbf{x}, \mathbf{w})) &= E_{c \sim p(c | \mathbf{x}, \mathbf{w})} [-\log p(c | \mathbf{x}, \mathbf{w})] \\ &\approx \frac{1}{N} \sum_{n=1}^N -\log p(c^n | \mathbf{x}, \mathbf{w}), \quad c^n \sim p(c | \mathbf{x}, \mathbf{w}). \end{aligned} \quad (6)$$

Each of these uncertainty measures based on the logarithmic score considers the distribution over all possible output sequences  $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ , which is defined over the entire set of possible output sequences  $\mathcal{Y}_T$ . To approximate this distribution, it requires sampling output sequences from  $\mathcal{Y}_T$ . This requires generating multiple output sequences, which is computationally expensive. In the following, we eliminate this requirement by considering an alternative proper scoring rule.

## 2.2 NEW UNCERTAINTY MEASURES IN NLG BASED ON ZERO-ONE SCORE

Next, we introduce measures based on the zero-one score, which has not yet been considered as a proper scoring rule for deriving uncertainty measures in NLG. The zero-one score considers the predictive distribution for the most likely output sequence:

$$S_{0-1}(p(\mathbf{y} \mid \mathbf{x}, \cdot), \mathbf{y}') = \begin{cases} 1 - p(\mathbf{y} = \mathbf{y}' \mid \mathbf{x}, \cdot) & \text{if } \mathbf{y}' = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}, \cdot), \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

Using the zero-one score in Eq. (2) results in the total uncertainty being the expected confidence of the given language model about the most likely output sequences generated by all language models according to their posterior probability  $p(\mathbf{w} \mid \mathcal{D})$ :

$$\begin{aligned} & \underbrace{E_{\tilde{\mathbf{w}} \sim p(\tilde{\mathbf{w}} \mid \mathcal{D})} [1 - p(\mathbf{y} = \tilde{\mathbf{y}}^* \mid \mathbf{x}, \mathbf{w})]}_{\text{total}} \\ &= \underbrace{1 - p(\mathbf{y} = \mathbf{y}^* \mid \mathbf{x}, \mathbf{w})}_{\text{aleatoric}} + \underbrace{p(\mathbf{y} = \mathbf{y}^* \mid \mathbf{x}, \mathbf{w}) - E_{\tilde{\mathbf{w}} \sim p(\tilde{\mathbf{w}} \mid \mathcal{D})} [p(\mathbf{y} = \tilde{\mathbf{y}}^* \mid \mathbf{x}, \mathbf{w})]}_{\text{epistemic}}, \end{aligned} \quad (8)$$

with  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$  and  $\tilde{\mathbf{y}}^* = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}, \tilde{\mathbf{w}})$ . Similar to Eq. (4), the epistemic uncertainty is a posterior expectation that remains challenging to estimate. However, we again focus on the aleatoric uncertainty, which solely considers the likelihood of the most likely output sequence under the given language model.

While aleatoric uncertainty derived from the logarithmic score requires approximating the entire output sequence distribution by sampling multiple sequences (as seen in Eq. (5) and Eq. (6)), the aleatoric uncertainty based on the zero-one score (see Eq. (8)) requires approximating the most likely output sequence under the given language model. This distinction is crucial, as approximating the most likely output sequence aligns directly with standard inference techniques widely used in language models, such as greedy decoding, beam search (Sutskever et al., 2014), top-k sampling, or nucleus sampling (Holtzman et al., 2020). For numerical stability, we consider the negative log-likelihood of the most likely output sequence that is proportional to the measure of aleatoric uncertainty in Eq. (8). We propose to estimate this quantity using the greedily decoded output sequence as an efficient and effective measure of aleatoric uncertainty:

$$\text{NLL} := - \sum_{t=1}^T \log \left( \max_{y_t} p(y_t \mid \mathbf{x}, \mathbf{y}_{<t}, \mathbf{w}) \right) \approx - \log p(\mathbf{y} = \mathbf{y}^* \mid \mathbf{x}, \mathbf{w}) \quad (9)$$

**Discussion.** Our proposed uncertainty measure challenges the prevailing reliance on multi-sequence sampling and semantic clustering for uncertainty estimation in NLG. By solely relying on the output sequences generated with greedy decoding, our approach significantly reduces computational overhead while maintaining theoretical rigor through its foundation in proper scoring rules. While uncertainty measures based on the logarithmic score could theoretically excel if the full distribution over output sequences  $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$  were accessible – as in standard classification tasks – this distribution is intractable for NLG tasks due to their sequential nature. As a result, sampling-based methods often yield crude approximations, constrained by computational limits and sampling variability. In contrast, our uncertainty measure, based on the zero-one score, offers a more rigorous alternative while eliminating the need for extensive sampling. In Sec. 4, we demonstrate that using our measure of uncertainty yields performance that is superior to or at least on par with uncertainty measures based on the logarithmic score. This makes our method more practical for large-scale applications.

### 3 RELATED WORK

In the previous section, we discussed uncertainty estimation methods based on the logarithmic score. Beyond these, there is a body of work that extends the concept of Semantic Entropy (Kuhn et al., 2023; Farquhar et al., 2024), for instance by either improving the semantic clustering (Nikitin et al., 2024; Qiu & Miiikkulainen, 2024), improving the sampling of output sequences (Aichberger et al., 2024), or directly approximating the measure from hidden states of the language model (Kossen et al., 2024). Also, there is a body of work that builds upon the concept of Predictive Entropy (Malinin & Gales, 2021), 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., 2023; Bakman et al., 2024).

There is also work on uncertainty estimation in NLG that is not grounded in proper scoring rules. For instance, several approaches leverage the language model itself to directly predict uncertainty, whether 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) employ cross-examination, where one language model generates an output sequence and another model acts as an examiner to assess uncertainty. Zhou et al. (2023) explore the behavior of language models when expressing their uncertainty, providing insights into how models articulate confidence in their predictions. Also, Manakul et al. (2023) propose using sampled output sequences as input for another language model to assess uncertainty, offering a unique perspective on sequence evaluation. Additionally, Xiao et al. (2022) provide an empirical analysis of how factors such as model architecture and training data influence uncertainty estimates. Finally, conformal prediction (Quach et al., 2023) offers another approach by calibrating a stopping rule for output sequence generation, providing a statistical framework for uncertainty estimation.

### 4 EXPERIMENTS

We aligned the evaluation of uncertainty estimation methods with related work by focusing on free-form question-answering tasks (Kuhn et al., 2023; Duan et al., 2023; Bakman et al., 2024; Nikitin et al., 2024; Aichberger et al., 2024; Kossen et al., 2024). While Farquhar et al. (2024) additionally concerns experiments with paragraph-length generations, their approach involves breaking down the entire paragraph into factual claims and reconstructing corresponding questions. Since the performance is expected to correlate with the performance on free-form question answering, we decided to focus specifically on free-form question answering tasks for a more direct assessment and less ambiguity in the evaluation.

**Datasets.** We evaluated uncertainty estimation methods on three different datasets. We used the over 3,000 test instances from *TriviaQA* (Joshi et al., 2017) concerning trivia questions, the over 300 test instances from *SVAMP* (Patel et al., 2021) concerning elementary-level math problems, and the over 3,600 test instances from *NQ-Open* (Lee et al., 2019) to assess natural questions aggregated from Google Search. Each dataset was utilized for two distinct tasks: (1) generating concise answers in the form of short phrases, and (2) producing more detailed answers in the form of full sentences (Farquhar et al., 2024). The resulting six tasks span a broad range of scenarios, ensuring a comprehensive evaluation of the uncertainty estimation methods.

**Models.** We conducted our evaluations on six distinct language models across different architectures, sizes, and training stages. Specifically, we used the Transformer-based model series *Llama-3.1* (Vaswani et al., 2017; Dubey et al., 2024) and the state-space model series *Falcon Mamba* (Gu & Dao, 2024; Zuo et al., 2024), representing two prominent language model paradigms. To assess the effect of training stage model scale on uncertainty estimation in NLG, we considered pre-trained (*PT*) and instruction-tuned (*IT*) language models with 7, 8, and 70 billion parameters, together covering a wide spectrum of model characteristics.

**Baselines.** We compare our method against the commonly used uncertainty measures based on the logarithmic score as of Eq. (5) and Eq. (6). These include Predictive Entropy (*PE*), length-normalized Predictive Entropy (*LN-PE*) (Malinin & 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. The discrete variant of Semantic

Table 1: Average AUROC across TriviaQA, SVAMP and NQ datasets, using uncertainty estimates of different measures to distinguish between correct and incorrect answers. Varying model architectures (*transformer*, *state-space*), model sizes (*7B*, *8B*, *70B*), and model stages (*PT*, *IT*) are considered for generating answers. The reference answer is generated using *greedy decoding*, either as a whole sentence (*long*) or a short phrase (*short*). The reference answer’s correctness is assessed by checking if the F1 score of the commonly used SQuAD metric exceeds 0.5 (*F1*) or if the LLM-as-a-judge considers it as correct (*LLM*). Predictive Entropy (*PE*), length-normalized Predictive Entropy (*LN-PE*), Semantic Entropy (*SE*), length-normalized Semantic Entropy (*LN-SE*), and discrete Semantic Entropy (*D-SE*) use 10 output sequences to assign an uncertainty estimate, each generated via multinomial sampling. NLL solely uses the reference answer to assign an uncertainty estimate.

		Uncertainty measure based score			Logarithmic				Zero-One	
Model	Gen.	Metric	PE	LN-PE	SE	LN-SE	D-SE	NLL		
Transformer	8B	short	F1	.776	.795	.775	.793	.804	.824	
		PT	short	LLM	.698	.714	.690	.706	.719	.726
			long	LLM	.562	.555	.545	.553	.600	.649
			short	F1	.772	.801	.805	.814	.806	.838
		IT	short	LLM	.676	.697	.704	.709	.694	.722
			long	LLM	.551	.548	.599	.601	.609	.615
	70B		short	F1	.775	.790	.793	.803	.791	.820
		PT	short	LLM	.693	.709	.718	.722	.715	.723
			long	LLM	.552	.534	.558	.569	.571	.649
			short	F1	.748	.781	.790	.799	.783	.792
		IT	short	LLM	.681	.698	.703	.709	.699	.699
			long	LLM	.555	.557	.568	.595	.600	.562
State-Space	7B	short	F1	.811	.815	.809	.822	.828	.843	
		PT	short	LLM	.705	.711	.701	.711	.716	.728
			long	LLM	.567	.597	.574	.611	.624	.612
			short	F1	.793	.814	.797	.816	.829	.838
		IT	short	LLM	.690	.701	.689	.699	.711	.719
			long	LLM	.588	.587	.597	.618	.629	.615

Entropy entirely disregards the output sequence likelihood and only considers the proportion of output sequences that belong to the same semantic cluster (Farquhar et al., 2024).

**Evaluation.** Effective uncertainty measures should accurately reflect the reliability of answers generated by the language model. Higher uncertainty more likely leads to incorrect generations. Thus, to evaluate the performance of an uncertainty estimator, we assess how well it correlates with the correctness of the language model’s answers; correct answers should be assigned a lower uncertainty estimator than incorrect answers. To determine whether an answer is correct, it has to be compared to the respective ground truth answer. To do so, we check if the F1 score of the commonly used SQuAD metric exceeds 0.5 (Rajpurkar et al., 2016). Although there are some limitations to using such a simple metric, it has relatively small errors in standard data sets and, therefore, remains widely used in practice. However, this metric is only applicable for short-phrase generations that align with the ground truth answer. Therefore, we additionally employ Llama-3.1 with 70 billion parameters (Dubey et al., 2024) as an LLM-as-a-judge to assess the correctness of both short-phrase and full-sentence generations. Subsequently, to measure the correlation between incorrectness of answers and the respective uncertainty estimates, we use the Area Under the Receiver Operating Characteristic (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. Overall, this evaluation process follows established methodologies for assessing the performance of uncertainty measures in NLG (Kuhn et al., 2023; Duan et al., 2023; Bakman et al., 2024; Farquhar et al., 2024; Nikitin et al., 2024; Aichberger et al., 2024; Kossen et al., 2024).

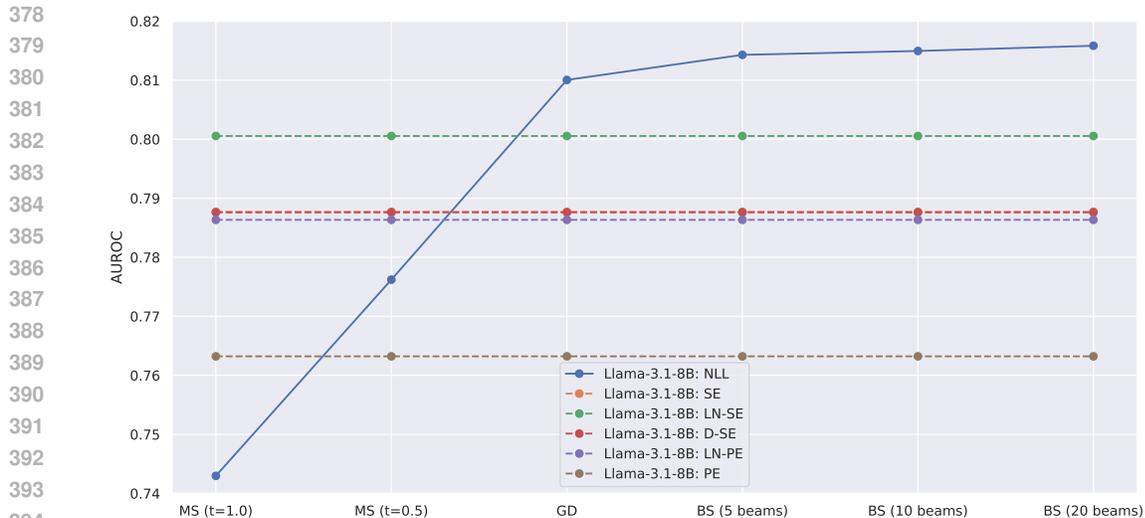


Figure 1: Average AUROC for the TriviaQA dataset, using the Llama-3.1-8B model to generate short phrase answers. The reference answer is generated using multinomial sampling (MS) with different temperature values ( $t$ ), Greedy decoding (GD), and beam search (BS) with a different number of beams.

**Analysis of results.** Tab. 1 summarizes the performance of uncertainty measures across six different language models and six different tasks. Our proposed measure (NLL) largely outperforms current state-of-the-art uncertainty measures, particularly in tasks that involve generating short phrases. This suggests that our measure is highly effective when focusing on the critical part of the output sequence that contains the actual answer to a question. In practical scenarios, the reliability of the specific answer is often more relevant than the uncertainty of the entire generated sentence. Thus, our measure provides targeted and computationally efficient uncertainty estimates, delivering enhanced performance where it is most critical, especially in real-world applications.

**Approximating the most likely output sequence.** Figure 1 illustrates the performance of our uncertainty measure when considering different inference techniques for generating answers. The reference answer, generated via beam search with a size of 20, is used to assess correctness, as it provides the best approximation of the most likely answer generated by the language model. Since the baselines are evaluated on output sequences generated using their optimal hyperparameter settings, their performance remains consistent. The results show that as the approximation to the most likely answer improves, so does the performance of our measure. However, while multinomial sampling significantly degrades the performance of our uncertainty measure, greedy decoding achieves performance comparable to more precise methods, such as beam search, reinforcing its validity as an effective approximation of the most likely output sequence.

Further experimental results and insights into the behavior of the uncertainty estimators can be found in Sec. A and Sec. B in the appendix.

## 5 CONCLUSION

We introduced a computationally efficient, theoretically grounded uncertainty measure, the negative log-likelihood of the most likely output sequence under a given language model. This measure is motivated by the general notion of proper scoring rules, providing a theoretically justified measure that is well aligned with the practical usage of LLMs. The experiments show that our measure performs extremely well with just a single generated output sequence, compared to previous measures that require multiple costly sequences to estimate the uncertainty. As a result, our approach represents a significant advance toward providing reliable uncertainty estimates that can be effectively applied at scale.

Although our proposed measure effectively captures uncertainty, it currently does not consider the semantics of the generated output sequence. Future work should investigate how it could be extended to also account for semantic meaning, **to further enrich the uncertainty estimator while preserving its computational efficiency.** Furthermore, all measures based on proper scoring rules depend on heuristics such as length normalization to deal with varying sequence lengths (Malinin & Gales, 2021; Duan et al., 2023; Bakman et al., 2024). **Investigating theoretically justified means to account for these varying generation characteristics is another promising direction for future work.** While there remain opportunities for refinement, **our proposed measure establishes a solid foundation for reliable and scalable uncertainty estimation in NLG.**

## ETHICS AND REPRODUCIBILITY STATEMENT

We acknowledge that language models can generate biased or harmful content if not properly managed. While our uncertainty estimation method enhances reliability, we encourage the responsible use of our approach in conjunction with bias mitigation and content moderation techniques.

We have made concerted efforts to ensure the reproducibility of our results. We report the raw average scores across held-out test datasets without standard error, as the distributional characteristics are more reflective of the models and datasets selected than the uncertainty estimation method itself. Theoretical derivations are provided in Sec. 2. All datasets are publicly available, and standard benchmarks are utilized to facilitate replication. The source code for reproducing all experiments will be made available upon publication.

## REFERENCES

- Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. Semantically diverse language generation for uncertainty estimation in language models. *arXiv*, 2406.04306, 2024.
- 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. *arXiv*, 2402.11756, 2024.
- Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. Crawling the internal knowledge-base of language models. In Andreas Vlachos and Isabelle Augenstein (eds.), *Findings of the Association for Computational Linguistics: EACL 2023*, pp. 1856–1869, Dubrovnik, Croatia, May 2023a. Association for Computational Linguistics.
- 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 (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12621–12640, Singapore, December 2023b. Association for Computational Linguistics.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Chenan Wang, Alex Zavalny, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the uncertainty estimation of large language models. *arXiv*, 2307.01379, 2023.
- 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. 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, Kamilè Lukošiušė, 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,

- 486 Nova DasSarma, Oliver Rausch, Robert Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu,  
487 Saurav Kadavath, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy  
488 Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann,  
489 Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Christopher Olah, Jack Clark,  
490 Samuel R. Bowman, and Jared Kaplan. The capacity for moral self-correction in large language  
491 models. *arXiv*, 2302.07459, 2023.
- 492 Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
493 *Journal of the American statistical Association*, 102(477):359–378, 2007.
- 494
- 495 Sebastian Gruber and Florian Buettner. Uncertainty estimates of predictions via a general bias-  
496 variance decomposition. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.),  
497 *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume  
498 206 of *Proceedings of Machine Learning Research*, pp. 11331–11354. PMLR, 25–27 Apr 2023.
- 499
- 500 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024.
- 501 Paul Hofman, Yusuf Sale, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty  
502 with proper scoring rules. *arXiv*, 2404.12215, 2024.
- 503
- 504 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text  
505 degeneration. In *International Conference on Learning Representations*, 2020.
- 506
- 507 Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for  
508 classification and preference learning. *arXiv*, 1112.5745, 2011.
- 509
- 510 Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning:  
511 An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.
- 512 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly  
513 supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meet-*  
514 *ing of the Association for Computational Linguistics*, Vancouver, Canada, July 2017. Association  
515 for Computational Linguistics.
- 516
- 517 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,  
518 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer  
519 El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bow-  
520 man, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna  
521 Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom  
522 Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Ka-  
523 plan. Language models (mostly) know what they know. *arXiv*, 2207.05221, 2022.
- 524
- 525 Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer  
526 vision? In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,  
527 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran  
Associates, Inc., 2017.
- 528
- 529 Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa Schut, Shreshth Malik, and Yarin Gal. Se-  
530 mantic entropy probes: Robust and cheap hallucination detection in llms. 2024.
- 531
- 532 Nikita Kotelevskii and Maxim Panov. Predictive uncertainty quantification via risk decompositions  
533 for strictly proper scoring rules. *arXiv*, 2402.10727, 2024.
- 534
- 535 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for  
536 uncertainty estimation in natural language generation. In *The Eleventh International Conference  
on Learning Representations*, 2023.
- 537
- 538 Meelis Kull and Peter Flach. Novel decompositions of proper scoring rules for classification: Score  
539 adjustment as precursor to calibration. In Annalisa Appice, Pedro Pereira Rodrigues, Vítor Santos  
Costa, Carlos Soares, João Gama, and Alípio Jorge (eds.), *Machine Learning and Knowledge  
Discovery in Databases*, pp. 68–85, Cham, 2015. Springer International Publishing.

- 540 Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym  
541 Korablyov, and Yoshua Bengio. DEUP: Direct epistemic uncertainty prediction. *Transactions on*  
542 *Machine Learning Research*, 2023. ISSN 2835-8856.
- 543  
544 Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised  
545 open domain question answering. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.),  
546 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp.  
547 6086–6096, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/  
548 v1/P19-1612.
- 549 Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in  
550 words. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856.
- 551 Andrey Malinin. *Uncertainty estimation in deep learning with application to spoken language*  
552 *assessment*. PhD thesis, University of Cambridge, 2019.
- 553  
554 Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In  
555 *International Conference on Learning Representations*, 2021.
- 556 Potsawee Manakul, Adian Liusie, and Mark Gales. SelfCheckGPT: Zero-resource black-box hal-  
557 lucination detection for generative large language models. In Houda Bouamor, Juan Pino, and  
558 Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Lan-*  
559 *guage Processing*, pp. 9004–9017, Singapore, December 2023. Association for Computational  
560 Linguistics.
- 561 Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational  
562 agents’ overconfidence through linguistic calibration. *Transactions of the Association for Com-*  
563 *putational Linguistics*, 10:857–872, 2022.
- 564  
565 Alexander Nikitin, Jannik Kossen, Yarín Gal, and Pekka Marttinen. Kernel language entropy: Fine-  
566 grained uncertainty quantification for llms from semantic similarities. 2024.
- 567  
568 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple  
569 math word problems? In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek  
570 Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou  
571 (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for*  
572 *Computational Linguistics: Human Language Technologies*, pp. 2080–2094, Online, June 2021.  
573 Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.168.
- 574  
575 Xin Qiu and Risto Miikkulainen. Semantic density: Uncertainty quantification in semantic space  
576 for large language models. 2024.
- 577  
578 Victor Quach, Adam Fisch, Tal Schuster, Adam Yala, Jae Ho Sohn, Tommi S. Jaakkola, and Regina  
579 Barzilay. Conformal language modeling. *arXiv*, 2306.10193, 2023.
- 580  
581 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
582 models are unsupervised multitask learners. 2018.
- 583  
584 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions  
585 for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), *Pro-*  
586 *ceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp.  
587 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi:  
588 10.18653/v1/D16-1264.
- 589  
590 Jie Ren, Yao Zhao, Tu Vu, Peter J. Liu, and Balaji Lakshminarayanan. Self-evaluation improves  
591 selective generation in large language models. *arXiv*, 2312.09300, 2023.
- 592  
593 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 (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 19446–19484. Curran Associates, Inc., 2023.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. On information-theoretic measures of predictive uncertainty, 2024.

- 594 Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks.  
595 *arXiv*, 1409.3215, 2014.  
596
- 597 Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea  
598 Finn, and Christopher Manning. Just ask for calibration: Strategies for eliciting calibrated confi-  
599 dence scores from language models fine-tuned with human feedback. In Houda Bouamor, Juan  
600 Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natu-  
601 ral Language Processing*, pp. 5433–5442, Singapore, December 2023. Association for Computa-  
602 tional Linguistics.
- 603 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
604 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,  
605 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy  
606 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,  
607 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel  
608 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,  
609 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,  
610 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,  
611 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh  
612 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen  
613 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,  
614 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models.  
615 *arXiv*, 2307.09288, 2023.
- 616 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
617 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg,  
618 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural  
619 Information Processing Systems*, volume 30, 2017.
- 620 Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-  
621 Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scale  
622 empirical analysis. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the  
623 Association for Computational Linguistics: EMNLP 2022*, pp. 7273–7284, Abu Dhabi, United  
624 Arab Emirates, December 2022. Association for Computational Linguistics.
- 625 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-  
626 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt  
627 Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer.  
628 Opt: Open pre-trained transformer language models. *arXiv*, 2205.01068, 2022.  
629
- 630 Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. Navigating the grey area: How expressions  
631 of uncertainty and overconfidence affect language models. In Houda Bouamor, Juan Pino, and  
632 Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Lan-  
633 guage Processing*, pp. 5506–5524, Singapore, December 2023. Association for Computational  
634 Linguistics.
- 635 Jingwei Zuo, Maksim Velikanov, Dhia Eddine Rhaiem, Ilyas Chahed, Younes Belkada, Guillaume  
636 Kunsch, and Hakim Hacid. Falcon mamba: The first competitive attention-free 7b language  
637 model. 2024.  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## A COMPARISON OF ESTIMATORS

In this section we want to empirically investigate the performance of estimators for the predictive entropy  $H(p(\mathbf{y} | \mathbf{x}))$  (Eq. (5)) and the maximum likelihood  $1 - \max_{\mathbf{y}} p(\mathbf{y} | \mathbf{x})$  (Eq. (8)). Therefore, we consider a synthetic experiment with the following setup. We are given a space of possible outcomes  $\mathcal{V}$  with  $|\mathcal{V}| = \{10, 100\}$ . 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  $p(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$  are not represented by a neural network, but randomly sampled (but fixed per run) 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 language models that follow a power law. For  $|\mathcal{V}| = 10$  we have  $\alpha_{1,2} = 10$  and  $\alpha_{3-10} = 0.2$ . 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 the leftmost subfigures in Fig. 2 and Fig. 3.

The experiments investigate the quality of the estimators depending on the number of samples  $\{\mathbf{y}_n\}_{n=1}^N$ . This is possible because it is possible to calculate the ground truth values for both the entropy and the maximum likelihood sequence for this small synthetic example by exhaustive enumeration. We average over 1,000 runs, meaning that the predictive distributions are redrawn according to the respective Dirichlet distribution. This corresponds to evaluating uncertainty for different input sequences  $\mathbf{x}$  for language models.

**Entropy estimation.** The results are shown in Fig. 2. We observe that the variance of estimators increases for larger vocabulary sizes  $|\mathcal{V}|$  and sequence lengths  $T$ . Furthermore, lower temperatures decrease the variance of the estimator at the cost of introducing bias.

**Maximum Likelihood.** The results are shown in Fig. 3. We observe that low-temperature multinomial sampling and beam search find the maximum log-likelihood with a low number of samples with high probability. Greedy decoding (beam size = 1) finds the maximum for all settings except the hardest ( $|\mathcal{V}| = 100, T = 4$ ), where it takes a beam size of 2 to find it.

## B DETAILED RESULTS

In this section, we provide detailed results to complement the main results presented in Tab. 1.

The main results used greedy decoding (beam search of size 1) to estimate the maximum likelihood (zero-one score based measure) and 10 samples to estimate entropies (logarithmic score based measures). For each dataset, we performed a hyperparameter search on held-out instances to determine the best performing temperature  $t \in \{0.5, 1.0, 1.5\}$  for sampling output sequences used for the logarithmic score based measures.

We look into how much the maximum likelihood benefits from additional samples by increasing the beam with to 5. The results are given in Tab. 2, showing that our measure continues to improve for a larger number of beams, thus better estimates of the maximum likelihood sequence. Furthermore, we provide detailed results for individual datasets in Tab. 3, complimenting the results presented in the main paper (c.f. Tab. 1).

The AUROC is considered as a primary performance measure throughout the paper. We additionally consider the average rejection accuracy, i.e. the accuracy of model predictions when allowing to reject a certain budget of predictions based on the uncertainty estimate. Thus, predictions are only evaluated for the 80% most certain predictions. Results are given in Tab. 4, again with greedy decoding for our measure based on the zero-one score. The results show, that our measure is very competitive across all settings, despite its simplicity and efficiency.

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

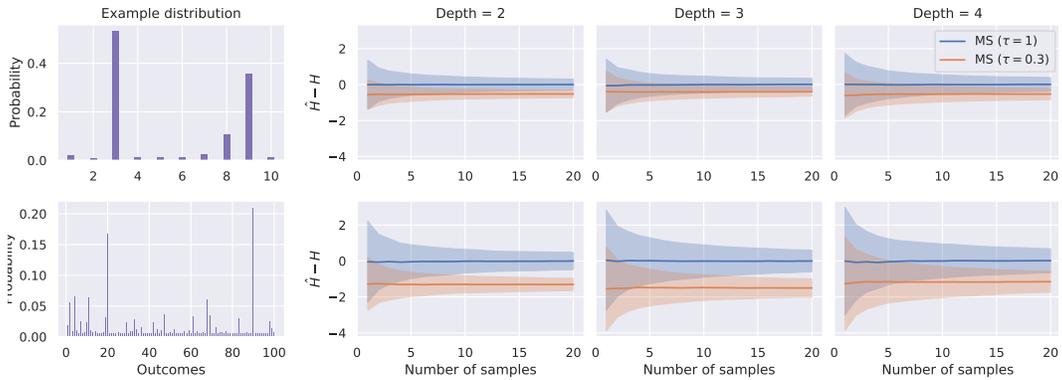


Figure 2: **Estimator of Predictive Entropy.** Results for different vocabulary sizes (rows) and sequence lengths (columns). The two leftmost subfigures show exemplary predictive distributions  $p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$ . We estimate the entropy using  $N$  samples by means of Eq. (5). Lines denote the average over runs, while shades denote one standard deviation. We compare multinomial sampling (MS) for two commonly used temperatures. The experiments show that temperature decreases variance but introduces bias.

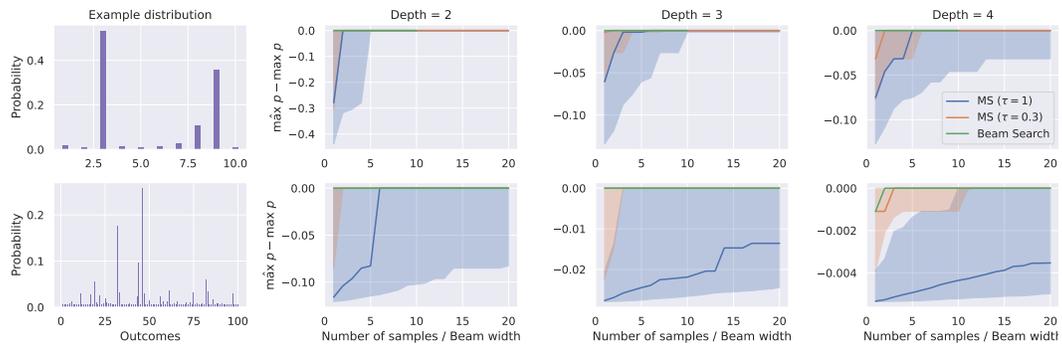


Figure 3: **Estimator of maximum likelihood.** Results for different vocabulary sizes (rows) and sequence lengths (columns). The two leftmost subfigures show exemplary predictive distributions  $p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$ . We estimate the maximum likelihood using the maximum over  $N$  sampled obtained by beam search or multinomial sampling (MS) with different temperatures. Lines denote the median, shades signify the possible values between the 5 and 95 percent quantile. Even with a very low number of samples, low-temperature multinomial sampling (MS) and beam search are able to find the maximum with high probability.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 2: **Average AUROC** across TriviaQA, SVAMP, and NQ datasets, using uncertainty estimates of different measures as a score to distinguish between correct and incorrect answers. Varying model architectures (*transformer*, *state-space*), model sizes (*7B*, *8B*, *70B*), and model stages (*PT*, *IT*) are considered for generating answers. The reference answer is generated using **beam search with 5 beams**, either as a whole sentence (*long*) or a short phrase (*short*). The correctness of the reference answer is assessed by checking if the F1 score of the commonly used SQuAD metric exceeds 0.5 (*F1*) or the Llama-3.1-70B model considers it as correct (*LLM*). Predictive Entropy (*PE*), length-normalized Predictive Entropy (*LN-PE*), Semantic Entropy (*SE*), length-normalized Semantic Entropy (*LN-SE*), and discrete Semantic Entropy (*D-SE*) use 10 output sequences to assign an uncertainty estimate, each generated via multinomial sampling. NLL solely uses the reference answer to assign an uncertainty estimate.

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

Table 3: **Average AUROC** of individual datasets, using uncertainty estimates of different measures as a score to distinguish between correct and incorrect answers.

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

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Table 4: **Average Rejection Accuracy (80%)** across TriviaQA, SVAMP and NQ 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, with the correctness being assessed by checking if the F1 score of the commonly used SQuAD metric exceeds 0.5 (*F1*), the pre-trained Llama-3.1-70B model considers it as correct (*LLM*), or the instruction-tuned Llama-3.1-70B-Instruct model considers it as correct (*LLM-Instruct*).

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