

Textual Bayes: Quantifying Uncertainty in LLM-Based Systems

Brendan Leigh Ross* Noël Vouitsis* Atiyeh Ashari Ghomi
 Rasa Hosseinzadeh Ji Xin Zhaoyan Liu Yi Sui Shiyi Hou
 Kin Kwan Leung Gabriel Loaiza-Ganem Jesse C. Cresswell
 {brendan, noel, atiyeh, rasa, ji, zhaoyan,
 amy, gloria, kk, gabriel, jesse}@layer6.ai
 Layer 6 AI, Toronto, Canada

Abstract

Although large language models (LLMs) are becoming increasingly capable of solving challenging real-world tasks, accurately quantifying their uncertainty remains a critical open problem—one that limits their applicability in high-stakes domains. This challenge is further compounded by the closed-source, black-box nature of many state-of-the-art LLMs. Moreover, LLM-based systems can be highly sensitive to the prompts that bind them together, which often require significant manual tuning (i.e., prompt engineering). In this work, we address these challenges by viewing LLM-based systems through a Bayesian lens. We interpret prompts as textual parameters in a statistical model, allowing us to use a small training dataset to perform Bayesian inference over these prompts. This novel perspective enables principled uncertainty quantification over both the model’s textual parameters and its downstream predictions, while also incorporating prior beliefs about these parameters expressed in free-form text. To perform Bayesian inference—a difficult problem even for well-studied data modalities—we introduce Metropolis-Hastings through LLM Proposals (MHLP), a novel Markov chain Monte Carlo (MCMC) algorithm that combines prompt optimization techniques with standard MCMC methods. MHLP is a turnkey modification to existing LLM pipelines, including those that rely exclusively on closed-source models. Empirically, we demonstrate that our method yields improvements in both predictive accuracy and uncertainty quantification (UQ) on a range of LLM benchmarks and UQ tasks. More broadly, our work demonstrates a viable path for incorporating methods from the rich Bayesian literature into the era of LLMs, paving the way for more reliable and calibrated LLM-based systems.

1 Introduction

Large language models (LLMs) have become increasingly embedded in our daily lives, with growing adoption across domains such as customer support [7], code generation [58, 8], scientific research [6, 52, 70], and creative writing [18]. As their capabilities continue to advance, there is also mounting interest in deploying them in agentic systems, wherein they perform tasks autonomously on behalf of users [67, 68]. Despite their rapid proliferation, however, trust in LLMs remains limited, largely due to their propensity to generate hallucinated content [38, 69] and their susceptibility to adversarial attacks and jailbreaking [61, 78, 71]. LLM-based systems therefore face important vulnerabilities that must be addressed to make them safe and trustworthy, especially when used in high-stakes domains such as finance and medicine. A key challenge towards mitigating these risks is to reliably quantify

*Equal contribution.

the uncertainty of LLM-based systems. Accurate measures of uncertainty ensure that LLM-based systems can abstain from providing false information, defer to human experts when appropriate, or selectively augment their context with subroutines based on retrieval or reasoning [32, 62].

Despite recent progress, UQ for LLMs is far from solved and no consensus exists over exactly what should be quantified [30, 60, 73]. In this work, we propose to better quantify uncertainty in LLM-based systems by viewing them through a Bayesian lens. In light of a model, observational data, and one’s prior beliefs, Bayesian inference uses Bayes’ rule to compute the distribution of possible model parameters. Commonly applied in classical statistics and deep learning, Bayesian inference is a principled and mathematically grounded approach to UQ [4]. Bayesian techniques have led to high-profile successes in methods like variational autoencoders [28] and Bayesian neural networks [5]. As with their 20th century forebears (e.g., [13, 51]), these methods estimate uncertainty over high-dimensional continuous variables. We bring Bayesian methods into the age of LLMs. In LLM-based systems, the main variables of interest are prompts, since LLMs themselves are often black boxes that can only be accessed via an API. By treating prompts as textual parameters in a statistical model, as illustrated in Figure 1, we can use Bayesian inference to estimate distributions over their values. These distributions rigorously quantify our uncertainty about the models themselves. Furthermore, they can be integrated into uncertainty estimates on the system’s downstream outputs via easy-to-compute Monte Carlo estimates. To the best of our knowledge, we are the first both to quantify the uncertainty associated with prompts in LLM-based systems, and to perform Bayesian inference for free-form textual variables.

Adapting Bayesian methods to text has its challenges and advantages. On the one hand, textual variables are discrete, making it difficult to apply traditional Bayesian deep learning techniques such as gradient-based Markov chain Monte Carlo (MCMC) [42, 64] or variational inference [51, 5]. We address this obstacle with a novel text-based MCMC method: *Metropolis-Hastings through LLM Proposals* (MHLP). On the other hand, textual variables are better suited conceptually to Bayesian modelling than high-dimensional continuous variables such as the weights of a deep neural network. Bayesian inference famously requires the specification of one’s prior beliefs about a variable; textual variables are more amenable to human priors than neural network weights, and as we show, prior beliefs can be readily incorporated into LLM-based systems as free-form text.

To advance and justify our proposed method, this work contains the following contributions:

1. We take a novel perspective on LLM-based systems in which prompts are viewed as Bayesian textual parameters θ in a model $p(y \mid x, \theta)$. We show how this formulation leads to a natural and principled way to incorporate prior beliefs about θ while quantifying the uncertainty inherent in the modelling process.
2. To implement our Bayesian approach, we propose *Metropolis-Hastings through LLM Proposals* (MHLP), an MCMC algorithm to sample from intractable distributions over textual variables. MHLP has broad potential applications even beyond Bayesian inference.
3. We propose a novel metric of model calibration, *semantic expected calibration error*, for quantifying calibration, a form of UQ, on free-form textual outputs.
4. We systematically evaluate our method through standard LLM benchmarks and baselines, showing that it improves performance while providing state-of-the-art UQ over model outputs.

2 Background and Terminology

2.1 LLM-Based Systems

The central object in this work is the LLM-based system. The most common LLM-based system is one consisting of a single input x (e.g., a question), a prompt θ (e.g., a system message defining instructions for the model’s behaviour, shared across all x), and an output y (e.g., the model’s predicted answer), which we denote

$$y = \mathbf{LLM}(x; \theta). \quad (1)$$

We allow $\mathbf{LLM}(x; \theta)$ to be any open- or closed-source model that we view as a *random* function of x and θ , whose randomness depends on the underlying LLM sampling strategy (e.g., greedy, temperature, nucleus). In general throughout the work, capitalized, boldface function names will indicate random functions comprising one or more LLM calls.

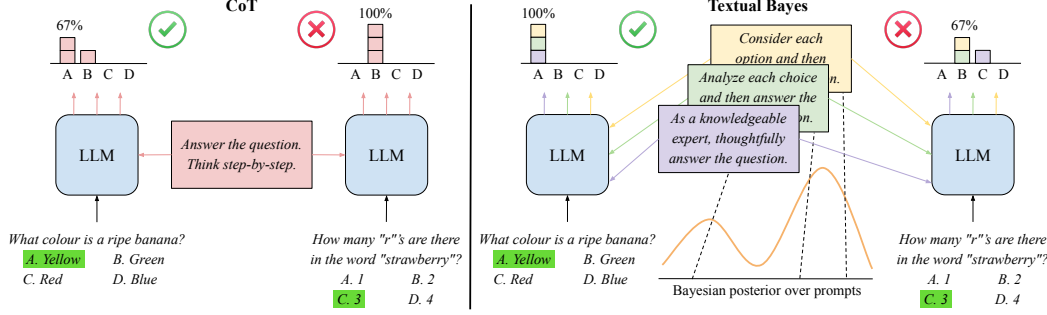


Figure 1: In chain-of-thought (CoT) prompting (left), answers are generated by an LLM using a single fixed prompt; this frequentist approach does not account for uncertainty about how the model should be prompted, causing potential issues such as overconfidence on incorrect answers. In our Textual Bayes approach (right), we sample prompts from our Bayesian posterior and then use each prompt to generate answers from the LLM; this allows for principled uncertainty quantification over both the prompts themselves and the resulting generated answers for improved model calibration.

LLM-based systems can be more complex than single-prompt models. Many recent works have proposed to group LLM calls of arbitrary count and complexity into pipelines parameterized by the prompts used at each step [27, 77, 74, 9, 22]. For example, Self-Refine [36] iterates on an initial LLM output by alternating between an LLM call providing feedback and one incorporating the feedback into refinement. Fully agentic systems integrate many LLM calls in sequence or parallel along with tool use to arrive at a final output. In full generality, we can describe a forward pass through an LLM-based system as

$$y = \mathbf{LBS}(x; \theta), \quad (2)$$

where $\mathbf{LBS}(\cdot; \theta)$ can be described as a directed acyclic graph with k edges in which each edge e_i corresponds to an LLM call $\mathbf{LLM}(\cdot; \theta_i)$ parameterized by a prompt θ_i and where we denote the combination of all prompts in the system as $\theta = (\theta_1, \dots, \theta_k)$. Since each LLM call in the system is potentially random, y is a random function of x parameterized by $\theta = (\theta_1, \dots, \theta_k)$. The LLM-based system thus forms a statistical model for y whose density we express as $p(y | x, \theta)$, where sampling $y \sim p(y | x, \theta)$ is equivalent to computing $y = \mathbf{LBS}(x; \theta)$.

Unlike a linear regressor or neural network where θ denotes continuous model parameters, for an LLM-based system θ denotes *textual* parameters. From the statistical modelling perspective, a natural next step is to find the optimal value of θ . For example, given an i.i.d. dataset $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$, one might want to perform maximum-likelihood:

$$\theta^* = \arg \max_{\theta} p(\mathcal{D} | \theta) = \arg \max_{\theta} \prod_{i=1}^n p(y_i | x_i, \theta). \quad (3)$$

The discrete nature of textual parameters prevents us from directly applying maximum-likelihood estimation to LLM-based systems via gradient-based algorithms. Prompt engineering can be understood as approximating θ^* by having a human propose candidate θ until adequate performance on a small dataset is reached. However, this manual process lacks rigour, is lengthy and tedious, and thus does not scale well.

Past works have proposed heuristic approaches to automatically optimize prompts θ in LLM-based systems [76, 27, 77, 9]. Here, we focus on iterative prompt optimization methods, which we can express mathematically as a stochastic update function **UPDATE**. Prompt optimization proceeds by iteratively applying **UPDATE** to an initial prompt $\theta^{(0)}$:

$$\theta^{(t)} = \mathbf{UPDATE}(\theta^{(t-1)}). \quad (4)$$

For simplicity, we assume that **UPDATE** is Markovian; i.e., not a function of θ values from earlier than $t - 1$. **UPDATE**, which consists of one or more LLM calls, is itself an LLM-based system.

One particularly relevant prompt optimization method is TextGrad [74]. The TextGrad framework conceptualizes constructive feedback on prompts as *textual gradients* and proposes a method for “backpropagating” feedback through an LLM-based system akin to backpropagation in neural networks. Although this framework is highly analogous to backpropagation of gradients for continuous variables, it does not formally optimize model likelihood.

2.2 Bayesian Inference

In this section, we briefly review Bayesian inference. For a more in-depth introduction, we refer the reader to [35]. Here, we allow $p(y | x, \theta)$ to be any statistical model for some variable y given another variable x . From the Bayesian perspective, there is uncertainty about the true value of θ , and hence the point estimate θ^* given by maximum-likelihood may be an overly reductive way of summarizing a dataset \mathcal{D} . Bayesian statistics provides a formal way of capturing this uncertainty.

First, we encode our prior uncertainty (beliefs) about the true value of θ as a *prior distribution* $p(\theta)$. Then, having observed a dataset \mathcal{D} , we update our beliefs about θ using Bayes' rule as

$$p(\theta | \mathcal{D}) = \frac{p(\theta)p(\mathcal{D} | \theta)}{p(\mathcal{D})} = \frac{p(\theta) \prod_i p(y_i | x_i, \theta)}{\sum_{\theta'} p(\theta') \prod_i p(y_i | x_i, \theta')}. \quad (5)$$

The *posterior distribution* $p(\theta | \mathcal{D})$ formally captures our uncertainty about θ in light of (i) our prior beliefs and (ii) the observed data. Given $p(\theta | \mathcal{D})$ and a new unobserved datapoint x_{new} , we can compute the *posterior predictive distribution* of y_{new} via

$$p(y_{\text{new}} | x_{\text{new}}, \mathcal{D}) = \sum_{\theta} p(y_{\text{new}} | x_{\text{new}}, \theta) p(\theta | \mathcal{D}) = \mathbb{E}_{\theta \sim p(\theta | \mathcal{D})} [p(y_{\text{new}} | x_{\text{new}}, \theta)]. \quad (6)$$

By Equation 6, we formalize predictive uncertainty in terms of uncertainty over θ . Since it is expressed as an expectation, we can estimate it via Monte Carlo sampling, assuming access to draws from $p(\theta | \mathcal{D})$. The posterior predictive has immediate practical value: estimates of its variability, such as variance or entropy, formally quantify uncertainty, and its value represents confidence in the prediction y_{new} .

The central challenge of Bayesian inference thus lies in sampling from the posterior $p(\theta | \mathcal{D})$. As $p(y | x, \theta)$ or $p(\theta)$ acquire even moderate complexity, sampling from $p(\theta | \mathcal{D})$ quickly becomes intractable. In deep learning, Bayesian inference requires approximations such as gradient-based MCMC [64], variational inference [5], or Laplace approximations [47]. We highlight that all of these approaches rely on the differentiability of $p(\theta)p(\mathcal{D} | \theta)$ with respect to θ , so none can be readily applied to the context where θ is a discrete prompt in an LLM-based system.

2.3 Markov Chain Monte Carlo and the Metropolis-Hastings Algorithm

In Bayesian statistics, MCMC algorithms are a common technique for tractably sampling from the posterior $p(\theta | \mathcal{D})$ when only its numerator in Equation 5 can be computed for any values of θ and \mathcal{D} . First, fix \mathcal{D} and let $g(\theta) = p(\theta)p(\mathcal{D} | \theta)$ be the numerator of Equation 5. Given an unnormalized density like $g(\theta)$, an MCMC algorithm is a general-purpose technique that specifies a Markov chain $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(t)}, \dots$ whose distribution converges to $\frac{g(\theta)}{\sum_{\theta'} g(\theta')} = p(\theta | \mathcal{D})$ as $t \rightarrow \infty$. In practice, by generating enough samples from the Markov chain, we can approximate sampling from $p(\theta | \mathcal{D})$ without needing to evaluate it.

The *Metropolis-Hastings* algorithm (MH) is a generic and broadly applicable form of MCMC (for a detailed introduction, see [49]). Starting with an initial sample $\theta^{(0)}$, MH iterates from sample $\theta^{(t-1)}$ to $\theta^{(t)}$ by generating a new *proposal* θ' from a pre-defined *proposal distribution* $q(\theta' | \theta)$ and then either accepting it (i.e., setting $\theta^{(t)} := \theta'$) or rejecting it (i.e., setting $\theta^{(t)} := \theta^{(t-1)}$) based on an acceptance probability γ (Algorithm 1).

In MH, the main “tunable hyperparameter”, the choice of proposal distribution $q(\theta' | \theta)$, is constrained only by very mild regularity conditions. However, the choice of q has a pronounced effect on the practicality of the algorithm, with poor choices of q (e.g., ones that perturb θ too mildly or too strongly

Algorithm 1: Metropolis-Hastings

Require: $\theta^{(0)}, q(\theta' | \theta), g(\theta)$;

for $t \leftarrow 1$ **to** T **do**

Sample proposal: $\theta' \sim q(\theta' | \theta^{(t-1)})$;

Compute acceptance probability:

$$\gamma = \min \left(1, \frac{g(\theta') q(\theta^{(t-1)} | \theta')}{g(\theta^{(t-1)}) q(\theta' | \theta^{(t-1)})} \right);$$

Sample random number: $u \sim \text{Uniform}(0, 1)$;

if $u < \gamma$ **then**

Accept: $\theta^{(t)} \leftarrow \theta'$;

else

Reject: $\theta^{(t)} \leftarrow \theta^{(t-1)}$;

return $\{\theta^{(t)}\}_{t=1}^T$;

at each step) taking an intractable amount of time to converge to the limiting distribution $p(\theta \mid \mathcal{D})$. The importance of q is such that some of the most popular MCMC algorithms (e.g., Langevin Monte Carlo and Hamiltonian Monte Carlo [13, 42]) are simply special cases of MH with highly specialized choices of q . Our method, MHLP, will also fall into this category, being specialized for textual parameters θ . The choice of q should be informed by any information available about the desired limiting distribution $p(\theta \mid \mathcal{D})$ [50]. Indeed, the optimal $q(\theta' \mid \theta)$ would be *equal* to the desired limiting distribution itself; if this were possible, of course, there would be no need to run MH in the first place. Nevertheless, we apply this intuition in Section 3 as we adapt MH to textual data.

3 Textual Bayes

In this section, we describe our method for Bayesian inference on LLM-based systems. We begin with the setup described in Section 2.1: an LLM-based system $\mathbf{LBS}(x; \theta)$ that gives rise to a statistical model $p(y \mid x, \theta)$, where x is the input, y is the output, and $\theta = (\theta_1, \dots, \theta_k)$ represents all the textual parameters involved in the system. We assume that $p(y \mid x, \theta)$ can be evaluated for any considered LLM-based system; see Appendix A for a discussion on this assumption and a description of how we can use open-source models as surrogates to estimate $p(y \mid x, \theta)$ and $p(\mathcal{D} \mid \theta)$ in practice. Our Bayesian inference algorithm will provide samples $\theta^{(1)}, \dots, \theta^{(m)} \sim p(\theta \mid \mathcal{D})$, which can in turn be used to quantify uncertainty over the system’s outputs as per Equation 6.

Textual priors To perform Bayesian inference, we must specify our prior beliefs about θ in the form of a distribution $p(\theta)$. Although θ lies in an infinite and semantically complex space of discrete text, humans are well equipped to reason and express their beliefs about textual variables. For example, a practitioner’s prior about a prompt θ_i might be that it should describe the purpose of the corresponding LLM call, guidelines for how to solve the task at hand, and the expected structure of the output. To exploit this knowledge, we codify our beliefs about each parameter θ_i as a free-form human-written string of textual constraints s_i , and provide it to an LLM to model the resulting parameter as

$$\theta_i = \mathbf{LLM}(s_i; \text{“Generate an LLM prompt satisfying the given constraints.”}). \quad (7)$$

For simplicity, we construct our prior $p(\theta) = \prod_{i=1}^k p(\theta_i)$ by assuming that all textual variables are independent, but this setup can be easily generalized by specifying joint constraints over multiple parameters θ_i and modelling them in a single LLM call.

Metropolis-Hastings through LLM Proposals Having constructed our prior $p(\theta)$, we now need an algorithm to sample from $p(\theta \mid \mathcal{D})$. A generally applicable MCMC method for text could have wide-ranging applications even beyond Bayesian inference. To this end, we propose Metropolis-Hastings through LLM Proposals (MHLP), a text-specific variant of MH.

At the heart of MHLP is our proposal distribution. We could in theory achieve the correct limiting distribution through almost any arbitrary choice of $q(\theta' \mid \theta)$, like randomly replacing letters or words in θ . But it is easy to see that such a proposal would rarely change θ semantically and never converge in practice. Instead, to generate useful proposals, we turn to LLMs. Analogously to how Langevin Monte Carlo uses gradient computation to exploit differentiable structure on $p(\theta \mid \mathcal{D})$, MHLP uses LLM calls to exploit linguistic structure on $p(\theta \mid \mathcal{D})$. Ideally, $q(\theta' \mid \theta)$ should be as similar to $p(\theta \mid \mathcal{D})$ as possible. By this standard, as per the relationship $p(\theta \mid \mathcal{D}) \propto p(\mathcal{D} \mid \theta)p(\theta)$, samples $\theta' \sim q(\theta' \mid \theta)$ should roughly satisfy the following criteria:

1. θ' should satisfy all the constraints embodied by the prior $p(\theta')$;
2. θ' should provide strong downstream performance on \mathcal{D} as measured by $p(\mathcal{D} \mid \theta')$.

We take inspiration from the prompt optimization methods discussed in Section 2.1 and use suggestions from LLMs to propose values of θ' that implement these guidelines. The observation underpinning MHLP is that iterative prompt optimization methods can be used to propose high-quality candidates θ' . Here we recall our formalization of prompt optimization as an iterated stochastic update function **UPDATE** (Equation 4), and sample from $q(\theta' \mid \theta)$ by computing $\theta' = \mathbf{UPDATE}(\theta)$.² Note that since **UPDATE** is itself an LLM-based system, evaluation of $q(\theta' \mid \theta)$ also follows from

²Some prompt optimization methods, such as the momentum variant of TextGrad [74], make updates based on a history of multiple past θ values. MHLP can take advantage of such methods by running multiple steps of the optimizer per accept/reject decision, akin to Hamiltonian Monte Carlo.

the approach discussed in Appendix A. Although MHLP is agnostic to the underlying prompt optimization method, we use TextGrad [74] in our implementation. By analogy to numerical losses in standard gradient-based optimization, TextGrad optimizes objectives described in natural language. We can thus express criteria #1 and #2 as objectives in natural language and use TextGrad to propose improvements to θ based on these criteria. This choice of objectives is specific to Bayesian inference, but in MHLP they can be easily replaced or modified to suit any textual distribution, which broadens its potential impact. We demonstrate one such example in Section 4.2.

In summary, we define MHLP as the variant of MH (Algorithm 1) acting on textual parameters θ in which the proposal step $\theta' \sim q(\theta', \theta^{(t-1)})$ is defined as $\theta' = \text{UPDATE}(\theta^{(t-1)})$.

Approximations for performance and tractability The approach outlined thus far enables sampling from the *true* Bayesian posterior over prompts in an LLM-based system. However, as is common practice in Bayesian deep learning—where approximations are often necessary for performance and tractability (e.g., [5, 12])—we propose analogous modifications to Algorithm 1. A well-studied phenomenon in Bayesian deep learning is the *cold posterior effect* wherein sampling from the “tempered” posterior $p_\tau(\theta \mid \mathcal{D}) \propto p(\mathcal{D} \mid \theta)^{1/\tau} p(\theta)$ with $0 < \tau < 1$ often results in better empirical performance than the standard Bayesian posterior (i.e., $\tau = 1$) [66, 2, 16, 23, 44, 26, 41]. We extend MHLP to accommodate $\tau \neq 1$ by setting $g(\theta) = p(\mathcal{D} \mid \theta)^{1/\tau} p(\theta)$ in Algorithm 1.

For tractability, we employ batching to approximate marginalizing over internal stochasticity within the system $\text{LBS}(x; \theta)$; marginalizing over reasoning traces when computing $p(y \mid x, \theta)$; and estimating the likelihood $p(\mathcal{D} \mid \theta) = \prod_{i=1}^n p(y_i \mid x_i, \theta)$, which would otherwise require a forward pass over the entire dataset. For implementation details on these approximations, please refer to Appendix A.

4 Experiments

In this section, we empirically evaluate our proposed Textual Bayes method. Specifically, we aim to answer the following question: how does Bayesian inference on the prompts of an LLM-based system with our MHLP algorithm translate into the system’s downstream *predictive performance* and *uncertainty quantification* (UQ) abilities?

In Section 4.1, we demonstrate that our method outperforms comparable baselines in accuracy, calibration, and abstention capabilities on challenging LLM benchmarks. Then, in Section 4.2, we adapt Textual Bayes to factuality, a context with weaker supervision than traditional Bayesian inference, showing that it can effectively reduce hallucinations using the conformal factuality framework [40].

Implementation For each dataset, we use MHLP to generate samples $\theta^{(1)}, \dots, \theta^{(m)} \sim p(\theta \mid \mathcal{D})$ from a Markov chain of length T . To increase sample diversity we employ burn-in, in which a fixed number d of initial MCMC samples are discarded, and thinning, in which we take every h -th sample thereafter until m samples are obtained. Given a datapoint x_{new} , we sample values of y_{new} using Equation 6 and quantify uncertainty on the basis of these downstream outputs. For other implementation details such as settings of d , h , m , and other dataset-specific hyperparameters, see Appendix B.

4.1 Uncertainty quantification with Textual Bayes

Setup We consider the canonical LLM-based system consisting of a single LLM as defined by Equation 1. Hallucinations in such systems occur when a model responds confidently with incorrect or ungrounded information, an issue that can be combatted with calibration [25, 63]. Calibration refers to the quality of a model’s confidence score, or the probability it assigns to the correctness of its provided answer. In other words, calibration measures how well LLMs “know what they know”. Here, we test calibration in downstream responses resulting from Bayesian inference over the LLM’s prompt. We compute confidences by generating 10 responses from each system and measuring the frequency of each response. For MHLP, we initialize $\theta^{(0)}$ to be a generic chain-of-thought (CoT) [62] prompt of the form: “Answer the question. Think step-by-step.”.

Baselines We compare our method against four frequentist baselines. *Paraphrasing* and *System-Message* are two prompt perturbation methods proposed by [17]. These methods inject prompt

stochasticity by rephrasing the question or system prompt in question-answering context.³ To these we add two additional baselines: (i) CoT refers to sampling m predictions from $\mathbf{LBS}(x; \theta^{(0)})$, and (ii) TextGrad refers to first performing T steps of prompt optimization and then sampling m predictions from $\mathbf{LBS}(x; \theta^{(T)})$. To ensure a fair comparison against MHLP, we equate m across all methods and use the same value of T for both TextGrad and MHLP.

We highlight that, because our method is a means of *generating* a diverse answer set $y^{(1)}, \dots, y^{(m)}$, our baselines are methods designed specifically to do the same. This means we omit direct comparison to means of *summarizing* $y^{(1)}, \dots, y^{(m)}$ into uncertainty scores, such as semantic entropy [30] and other methods described in Section 5. Although these are also UQ methods, they are orthogonal to and can be straightforwardly combined with our approach. For direct comparison, all experiments in this section use confidence or semantic confidence (described below) as the means of summarizing the uncertainty in every set of answers.

Datasets We evaluate both predictive performance and model calibration on AIME 2024 [10], SimpleQA [63], and QASPER [11], representing question-answering tasks that are closed-form, free-form, and free-form with context, respectively. We randomly select and fix 100 samples from each of SimpleQA and QASPER for all experiments and use all 30 available samples from AIME 2024. Notably, QASPER includes questions explicitly marked as unanswerable for which no supporting context is provided. We leverage these instances to assess our method’s ability to generalize to out-of-distribution samples by evaluating cases where the model must detect insufficient information and abstain from answering. See Appendix B for further dataset details.

Table 1: Accuracy (%) across datasets

Method	AIME	SimpleQA	QASPER
Paraphrasing	12.6 \pm 0.7	43.7 \pm 0.5	43.7 \pm 1.3
System-Message	7.2 \pm 0.7	47.3 \pm 0.7	59.7 \pm 0.6
CoT	9.0 \pm 1.4	47.8 \pm 0.6	56.5 \pm 0.8
TextGrad	11.9 \pm 0.9	46.6 \pm 0.5	58.8 \pm 1.0
MHLP (Ours)	15.0 \pm 0.7	48.6 \pm 0.6	60.9 \pm 1.0

Table 2: ECE / SECE (%) across datasets

Method	AIME	SimpleQA	QASPER
Paraphrasing	21.1 \pm 0.8	18.7 \pm 0.7	28.5 \pm 1.1
System-Message	19.7 \pm 0.8	18.4 \pm 0.4	23.9 \pm 0.9
CoT	31.5 \pm 1.4	18.0 \pm 0.6	26.2 \pm 0.67
TextGrad	27.4 \pm 1.6	17.7 \pm 1.0	21.6 \pm 1.2
MHLP (Ours)	22.0 \pm 1.0	15.4 \pm 0.6	17.7 \pm 1.1

In Table 1, we report accuracy for all datasets using exact-match on closed-form datasets and an LLM judge [75] to assess semantic correctness on free-form datasets. In Table 2, we report the expected calibration error (ECE) as a measure of model calibration [45, 20]. Additionally, for QASPER, we estimate abstention ability on two types of unanswerable questions: questions with no context, and those with a random context. We use the same confidence scores used to estimate calibration as an abstention metric and compute the ROC AUC of this score when used as a classifier of answerability. Results are shown in Table 3. All results are averaged over 10 independent runs with standard errors to account for stochasticity.

Semantic ECE Standard ECE cannot be applied to open-ended tasks since it requires a confidence score, which is nontrivial to compute in general due to the variability of possible correct responses. To address this limitation, inspired by semantic entropy [30], we propose an extension of ECE based on semantic clustering. Our metric, semantic ECE (SECE), uses these clusters to estimate model confidence over free-form outputs. Specifically, for each input x_i , we sample m outputs: $y_i^{(1)}, \dots, y_i^{(m)} \sim p(y | x_i, \theta)$. We then query an LLM to group these outputs into semantic clusters. The empirical probability assigned by the model to each cluster is defined as the proportion of the generated samples in that cluster. The maximum of these probabilities is then taken as the model’s *semantic confidence* for input x_i . Finally, we use this value as the confidence for standard ECE computation, enabling estimation of model calibration for free-form outputs.

Table 3: Abstention ROC AUC (%)

Method	QASPER	
	No context	Random context
Paraphrasing	48.2 \pm 1.1	62.1 \pm 1.6
System-Message	76.6 \pm 1.7	69.9 \pm 1.3
CoT	75.6 \pm 1.1	67.4 \pm 0.9
TextGrad	66.6 \pm 2.1	67.4 \pm 0.9
MHLP (Ours)	77.9 \pm 1.2	71.7 \pm 0.9

Discussion Across tasks, MHLP is the only method to consistently outperform the rest. It only trails in calibration (ECE) on AIME, however its accuracy exceeds the two best-calibrated methods by a substantial margin. We hypothesize this outperformance is due to the high-posterior-valued samples of θ generated by MHLP; it effectively performs stochastic prompt optimization, incorporating training

³Our implementation has minor differences with the cited paper. For further details see Section B.2.1.

data performance into its accept/reject decisions. In contrast, TextGrad alone has no accept/reject scheme and thus “always accepts”, leading to the inclusion of potentially less useful changes to the initial prompt. For qualitative examples, see Appendix B.

4.2 Conformal factuality with MHL P

Background Conformal factuality [40] is a method for providing statistical guarantees on the correctness of LLM-generated answers to open-ended questions based on conformal prediction (CP) [55, 54]. Generally, CP techniques use a small set of n labeled datapoints to calibrate a prediction threshold. In conformal factuality, given a question x , an LLM generates an answer y which is broken into a set of distinct claims $\{c_1, \dots, c_\ell\}$. Each claim is assigned a factuality score $\mathbf{F}(c; \theta)$ —generated by an LLM-based system—with larger values indicating increased confidence that c is a factual claim. Then, after using CP to calibrate a threshold λ , claims with $\mathbf{F}(c; \theta) < \lambda$ are filtered out, such that only high-confidence claims are returned in the final answer \hat{y} . CP guarantees that \hat{y} contains only factual claims with high probability,

$$1 - \alpha \leq \mathbb{P}[c \text{ is factual } \forall c \in \hat{y}] \leq 1 - \alpha + \frac{1}{n+1}, \quad (8)$$

where the error rate α is user-defined. The quality of final answers can be gauged through the fraction of claims which are retained, since longer answers with more claims are more useful.⁴ Better calibrated expressions of confidence through $\mathbf{F}(c; \theta)$ improve claim retention, and since MHL P enables better calibration, we can use it to design a better factuality score.

Baseline (GPT-4 frequency scoring) The best performing option for $\mathbf{F}(c; \theta)$ from [40] is frequency scoring. Five alternative answers y_i are generated for the same question x from GPT-4 [1] using unit temperature and a manually crafted prompt θ to gauge self-consistency [57, 37]. For each claim in the original answer y , the number of times it appears across the y_i is used as the score.

Our method (MHL P frequency scoring) Like GPT-4 frequency scoring, our method estimates a claim’s importance based on its frequency across alternative generations. However, instead of generating with a single fixed prompt, we produce diverse alternatives by sampling different θ via MHL P with zero temperature. Notably, in the factuality context, ground truth outputs $\{y_1, \dots, y_n\}$ are unavailable, so the unnormalized posterior $p(\theta)p(\mathcal{D} \mid \theta)$ is unavailable. To surmount this obstacle, we replace the unnormalized probability mass with a surrogate

$$g(\theta) = \mathbb{E}_{p(y' \mid x, \theta)} \left[\frac{1}{|y'|} \sum_{c \in y'} \mathbf{F}(c; \theta) \right], \quad (9)$$

which we estimate stochastically when running Algorithm 1. One alternative answer y_i is generated per sampled prompt. The ability to sample from this surrogate distribution underscores MHL P’s versatility in situations beyond conventional Bayesian inference.

Dataset We use FactScore [39], which is widely adopted for evaluating factuality in open-ended LLM generations. Following [40], we focus on “person” entities from the biography generation subset and generate claims from the generated biographies using a fixed extraction method to ensure consistency across runs. We also follow the setup from [40] by using 50 samples for the calibration/test set, and perform 1000 random splits of calibration and test data for each value of α .

Implementation We initialize both scoring methods with the same prompt. For MHL P, we perform sampling using a separate set of 100 samples from FactScore and obtain five prompt samples. Since there is no ground-truth answer in the open-ended QA setting, factuality is determined by decomposing answers into claims (as in [40]) and annotating them using a GPT web search tool. We use GPT-4 for answer generation, and GPT-4o-mini for claim generation, factuality annotation, frequency scoring, and MHL P proposals. For more details on the implementation, see Appendix B.

Results First, we verify that both scoring methods achieve the target coverage from Equation 8: Figure 2a shows that empirical factuality remains within the conformal bounds across all values of α . Figure 2b compares the removal rate, with error bars showing the standard deviation of the average removal rate across the 1000 data splits. Our method consistently achieves lower removal, showing that MHL P scoring provides a better uncertainty estimation of the factuality of LLM outputs.⁵

⁴Filtering out all claims guarantees that \hat{y} does not contain false claims, but does not give a useful answer.

⁵The GPT-4 frequency scoring method shows slightly higher removal than reported in [40], likely due to our use of a stricter web search–based factuality annotator.

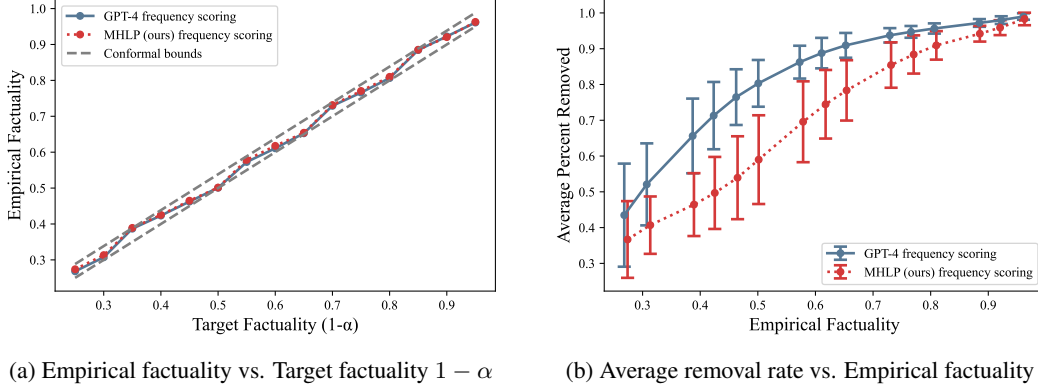


Figure 2: Comparison of conformal factuality for frequency scoring with a fixed prompt [40], and with prompts sampled through MHLP. (a) The empirical factuality achieved in practice is consistently within the bounds guaranteed by Equation 8. (b) MHLP achieves the same level of empirical factuality as frequency scoring but removes fewer claims, indicating better calibrated confidence.

5 Related Work

LLMs are applicable to a wide range of tasks and settings, which makes UQ inherently ambiguous—there is no single, well-defined quantity that UQ aims to approximate. Although our main setting of interest is where we have access to a pre-trained model and no fine-tuning is performed, we note that popular methods for UQ in deep learning, such as ensembles [31] and Laplace approximations [48, 29, 12], have been successfully ported over for UQ when fine-tuning LLMs [56, 72]. Within our setting of interest, some approaches estimate uncertainty by analyzing the variability in outputs generated by an LLM given the same input [30, 33, 19, 60, 46, 43], others do so by perturbing or modifying the input itself (e.g. by paraphrasing) [21, 17], and still others rely on directly asking the model to express its own confidence [25, 73]. Unlike all these methods, we aim to quantify the uncertainty associated with LLM prompts. Other methods for UQ within in-context learning tasks with LLMs have also leveraged Bayesian ideas [34, 24], but we highlight that these works differ greatly from ours in that they do not directly perform Bayesian inference over free-form text, and once again, they do not quantify uncertainty over prompts.

Lastly, we mention another line of work performing Metropolis-Hastings over text. Like ours, Faria et al. [15] use LLMs to construct a proposal distribution within the MH algorithm. We nonetheless highlight many differences with this work: their method is applied to machine translation and not to UQ, they do not perform Bayesian inference, and their proposal is completely different and does not rely on prompt optimization methods. Also, concurrently to our work, Faria and Smith [14] build on top of Faria et al. [15] by applying their proposal to a Bayesian formulation of the alignment problem wherein aligned model answers are sampled directly using MCMC. Investigating applications of our proposal to the alignment context or their proposal to UQ in LLM-based systems are potential directions for future work.

6 Conclusions, Limitations, and Future Work

In this work, we propose Textual Bayes for quantifying uncertainty in LLM-based systems. Our work represents a formalization of recent work conceptualizing LLM-based systems as models whose parameters are their prompts. Textual Bayes furthers this framework by performing Bayesian inference on these parameters, thus blending cutting-edge models with a formal statistical framework for uncertainty quantification. To implement this framework, we propose Metropolis-Hastings through LLM Proposals (MHLP), a novel MCMC algorithm for free-form text which finds applications in Bayesian inference and beyond. We test these frameworks on several uncertainty quantification benchmarks and find that they consistently improve the frontier of accuracy and calibration. We also show that MHLP can be adapted to a factuality-based objective, leading to more reliable factual claims as quantified by the setting of conformal factuality.

Although Textual Bayes and MHLP post strong performance against baselines, there remain limitations and avenues for improvement. First, MCMC is costly; despite consuming similar inference compute to leading baselines, Textual Bayes requires a one-time expensive application of MHLP. This cost might be addressed, for example, by further engineering the underlying prompt optimization method or training a small language model specifically for the task of generating proposals. Second, like many practical applications of Bayesian inference, our method involves several approximations which will inevitably cause samples to deviate from the true posterior. Third, our evaluations on free-form answering benchmarks require LLM-based judgements and LLM-based clustering. These techniques, though fairly applied across methods, are imperfect and a stronger evaluation signal might be obtained with improved fine-tuning, prompt engineering, or human evaluation. Lastly, we expect future work to find broader applications for MHLP beyond Bayesian inference. For example, we could use MHLP to modulate the *outputs* of LLM-based systems in accordance with unnormalized functions quantifying objectives such as alignment or safety.

References

- [1] Josh Achiam et al. GPT-4 Technical Report. *arXiv:2303.08774*, 2023.
- [2] Laurence Aitchison. A statistical theory of cold posteriors in deep neural networks. In *International Conference on Learning Representations*, 2021.
- [3] Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, et al. Llama-nemotron: Efficient reasoning models. *arXiv preprint arXiv:2505.00949*, 2025.
- [4] José M Bernardo and Adrian FM Smith. *Bayesian Theory*, volume 405. John Wiley & Sons, 2009.
- [5] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pages 1613–1622, 2015.
- [6] Daniil A Boiko, Robert MacKnight, and Gabe Gomes. Emergent autonomous scientific research capabilities of large language models. *arXiv:2304.05332*, 2023.
- [7] Rijul Chaturvedi and Sanjeev Verma. *Opportunities and Challenges of AI-Driven Customer Service*, pages 33–71. Springer International Publishing, 2023. ISBN 978-3-031-33898-4. doi: 10.1007/978-3-031-33898-4_3.
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv:2107.03374*, 2021.
- [9] Ching-An Cheng, Allen Nie, and Adith Swaminathan. Trace is the Next AutoDiff: Generative Optimization with Rich Feedback, Execution Traces, and LLMs. In *Advances in Neural Information Processing Systems*, volume 37, pages 71596–71642, 2024.
- [10] MAA Committees. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions.
- [11] Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. A dataset of information-seeking questions and answers anchored in research papers. *CoRR*, abs/2105.03011, 2021. URL <https://arxiv.org/abs/2105.03011>.

- [12] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux-effortless Bayesian deep learning. In *Advances in Neural Information Processing Systems*, 2021.
- [13] Simon Duane, A. D. Kennedy, B. J. Pendleton, and Duncan Roweth. Hybrid monte carlo. *Physics Letters B*, 195(2):216–222, 1987.
- [14] Gonalo Faria and Noah A Smith. Sample, don’t search: Rethinking test-time alignment for language models. *arXiv preprint arXiv:2504.03790*, 2025.
- [15] Gonalo Faria, Sweta Agrawal, Ant3nio Farinhas, Ricardo Rei, Jos3 de Souza, and Andr3 Martins. QUEST: Quality-aware metropolis-hastings sampling for machine translation. In *Advances in Neural Information Processing Systems*, 2024.
- [16] Vincent Fortuin, Adri3 Garriga-Alonso, Sebastian W. Ober, Florian Wenzel, Gunnar Ratsch, Richard E Turner, Mark van der Wilk, and Laurence Aitchison. Bayesian neural network priors revisited. In *International Conference on Learning Representations*, 2022.
- [17] Xiang Gao, Jiaxin Zhang, Lalla Mouatadid, and Kamalika Das. SPUQ: Perturbation-based uncertainty quantification for large language models. *arXiv preprint arXiv:2403.02509*, 2024.
- [18] Carlos G3mez-Rodr3guez and Paul Williams. A Confederacy of Models: a Comprehensive Evaluation of LLMs on Creative Writing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14504–14528, 2023.
- [19] Yashvir S Grewal, Edwin V Bonilla, and Thang D Bui. Improving uncertainty quantification in large language models via semantic embeddings. *arXiv:2410.22685*, 2024.
- [20] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1321–1330, 2017.
- [21] Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. Decomposing uncertainty for large language models through input clarification ensembling. In *International Conference on Machine Learning*, 2024.
- [22] Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *arXiv:2408.08435*, 2024.
- [23] Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Wilson. What Are Bayesian Neural Network Posteriors Really Like? In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 4629–4640, 2021.
- [24] Andrew Jesson, Nicolas Beltran Velez, Quentin Chu, Sweta Karlekar, Jannik Kossen, Yarin Gal, John P Cunningham, and David Blei. Estimating the hallucination rate of generative AI. In *Advances in Neural Information Processing Systems*, 2024.
- [25] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- [26] Sanyam Kapoor, Wesley J Maddox, Pavel Izmailov, and Andrew G Wilson. On uncertainty, tempering, and data augmentation in bayesian classification. In *Advances in Neural Information Processing Systems*, volume 35, 2022.
- [27] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. DSPy: Compiling declarative language model calls into state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations*, 2024.
- [28] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.

- [29] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being Bayesian, even just a bit, fixes overconfidence in relu networks. In *International Conference on Machine Learning*, 2020.
- [30] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *International Conference on Learning Representations*, 2023.
- [31] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, 2017.
- [32] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474, 2020.
- [33] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. In *Transactions on Machine Learning Research*, 2024.
- [34] Chen Ling, Xujiang Zhao, Xuchao Zhang, Wei Cheng, Yanchi Liu, Yiyu Sun, Mika Oishi, Takao Osaki, Katsushi Matsuda, Jie Ji, et al. Uncertainty quantification for in-context learning of large language models. *arXiv:2402.10189*, 2024.
- [35] David JC MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- [36] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-Refine: Iterative Refinement with Self-Feedback. In *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594, 2023.
- [37] Potsawee Manakul, Adian Liusie, and Mark Gales. SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.emnlp-main.557.
- [38] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, 2020.
- [39] Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.emnlp-main.741.
- [40] Christopher Mohri and Tatsunori Hashimoto. Language models with conformal factuality guarantees. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- [41] Seth Nabarro, Stoil Ganev, Adrià Garriga-Alonso, Vincent Fortuin, Mark van der Wilk, and Laurence Aitchison. Data augmentation in Bayesian neural networks and the cold posterior effect. In *Uncertainty in Artificial Intelligence*, pages 1434–1444. PMLR, 2022.
- [42] Radford M. Neal. *Bayesian Learning for Neural Networks*, volume 118 of *Lecture Notes in Statistics*. Springer, 1996. doi: 10.1007/978-1-4612-0745-0.
- [43] Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. Kernel language entropy: Fine-grained uncertainty quantification for LLMs from semantic similarities. In *Advances in Neural Information Processing Systems*, 2024.

- [44] Lorenzo Noci, Kevin Roth, Gregor Bachmann, Sebastian Nowozin, and Thomas Hofmann. Disentangling the roles of curation, data-augmentation and the prior in the cold posterior effect. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [45] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1), 2015. doi: 10.1609/aaai.v29i1.9602.
- [46] Xin Qiu and Risto Miikkulainen. Semantic density: Uncertainty quantification for large language models through confidence measurement in semantic space. In *Advances in Neural Information Processing Systems*, 2024.
- [47] Hippolyt Ritter, Aleksandar Botev, and David Barber. A Scalable Laplace Approximation for Neural Networks. In *International Conference on Learning Representations*, 2018.
- [48] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable Laplace approximation for neural networks. In *International Conference on Learning Representations*, 2018.
- [49] Christian P Robert. The Metropolis-Hastings algorithm. *arXiv:1504.01896*, 2015.
- [50] Jeffrey S Rosenthal. Optimal proposal distributions and adaptive MCMC. *Handbook of Markov Chain Monte Carlo*, 4(10.1201):93–111, 2011.
- [51] Lawrence K Saul, Tommi Jaakkola, and Michael I Jordan. Mean field theory for sigmoid belief networks. *Journal of Artificial Intelligence Research*, 4:61–76, 1996.
- [52] Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent Laboratory: Using LLM Agents as Research Assistants. *arXiv:2501.04227*, 2025.
- [53] Daniel Seita, Xinlei Pan, Haoyu Chen, and John Canny. An efficient minibatch acceptance test for metropolis-hastings. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 5359–5363, 2018.
- [54] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- [55] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, 2005.
- [56] Xi Wang, Laurence Aitchison, and Maja Rudolph. Lora ensembles for large language model fine-tuning. *arXiv:2310.00035*, 2023.
- [57] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [58] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8696–8708, 2021.
- [59] Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. Helpsteer2-preference: Complementing ratings with preferences, 2024. URL <https://arxiv.org/abs/2410.01257>.
- [60] Ziyu Wang and Chris Holmes. On subjective uncertainty quantification and calibration in natural language generation. *arXiv:2406.05213*, 2024.
- [61] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does LLM safety training fail? *Advances in Neural Information Processing Systems*, 36:80079–80110, 2023.

- [62] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837, 2022.
- [63] Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models. *arXiv preprint arXiv:2411.04368*, 2024.
- [64] Max Welling and Yee Whye Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. In *Proceedings of the 28th International Conference on Machine Learning*, pages 681–688, 2011. ISBN 978-1-4503-0619-5.
- [65] Bingbing Wen, Bill Howe, and Lucy Lu Wang. Characterizing llm abstention behavior in science qa with context perturbations, 2024. URL <https://arxiv.org/abs/2404.12452>.
- [66] Florian Wenzel, Kevin Roth, Bastiaan Veeling, Jakub Swiatkowski, Linh Tran, Stephan Mandt, Jasper Snoek, Tim Salimans, Rodolphe Jenatton, and Sebastian Nowozin. How good is the Bayes posterior in deep neural networks really? In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 10248–10259, 2020.
- [67] Michael Wooldridge and Nicholas R Jennings. Intelligent agents: Theory and practice. *The Knowledge Engineering Review*, 10(2):115–152, 1995.
- [68] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 68(2):121101, 2025.
- [69] Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models. *arXiv:2401.11817*, 2024.
- [70] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha. The AI Scientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search. *arXiv:2504.08066*, 2025.
- [71] Jun Yan, Vikas Yadav, Shiyang Li, Lichang Chen, Zheng Tang, Hai Wang, Vijay Srinivasan, Xiang Ren, and Hongxia Jin. Backdooring instruction-tuned large language models with virtual prompt injection. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6065–6086, 2024.
- [72] Adam X Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation for large language models. In *International Conference on Learning Representations*, 2024.
- [73] Daniel Yang, Yao-Hung Hubert Tsai, and Makoto Yamada. On verbalized confidence scores for LLMs. *arXiv:2412.14737*, 2024.
- [74] Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. Optimizing generative AI by backpropagating language model feedback. *Nature*, 639:609–616, 2025.
- [75] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- [76] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *The Eleventh International Conference on Learning Representations*, 2022.

- [77] Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. GPTSwarm: Language Agents as Optimizable Graphs. In *Forty-first International Conference on Machine Learning*, 2024.
- [78] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv:2307.15043*, 2023.

A Method Details

In general, MCMC can only be applied to Bayesian inference when the $g(\theta)$ is calculable, where $g(\theta)$ is defined by

$$g(\theta) = p(\theta)p(\mathcal{D} \mid \theta) = p(\theta) \prod_{i=1}^n p(y_i \mid x_i, \theta). \quad (10)$$

In our context, certain terms in this equation are intractable and we instead estimate them stochastically.

Mini-batch estimates of $\prod_{i=1}^n p(y_i \mid x_i, \theta)$ Computing the term $\prod_{i=1}^n p(y_i \mid x_i, \theta)$ requires evaluating $\mathbf{LBS}(x_i)$ for all $i \in \{1, \dots, n\}$. Training sets are often large enough to make running a full pass per MCMC step intractable. Instead, we use mini-batching. One way to mini-batch is to apply the MH adjustment of [53], which involves a different accept/reject step; however, to avoid complexity we use a simple stochastic estimate of $\prod_{i=1}^n p(y_i \mid x_i, \theta)$ instead. Using a batch size of b , we make the estimate

$$\prod_{i=1}^n p(y_i \mid x_i, \theta) \approx \prod_{j=1}^b p(y_{i_j} \mid x_{i_j}, \theta)^{\frac{n}{b}}, \quad (11)$$

which is an unbiased estimate in log-space:

$$\sum_{i=1}^n \log p(y_i \mid x_i, \theta) = n \mathbb{E}_{(x,y) \sim \text{Uniform}(\mathcal{D})} [\log p(y \mid x, \theta)] \approx \frac{n}{b} \sum_{j=1}^b \log p(y_{i_j} \mid x_{i_j}, \theta), \quad (12)$$

where $\{(x_{i_1}, y_{i_1}), \dots, (x_{i_b}, y_{i_b})\}$ is a mini-batch of training datapoints. In experiments, we use a batch size of $b = 1$. However, as discussed in Section 3, we also apply a temperature τ as in the cold posterior effect, making the final estimate equal to

$$\prod_{i=1}^n p(y_i \mid x_i, \theta)^{\frac{1}{\tau}} \approx \prod_{j=1}^b p(y_{i_j} \mid x_{i_j}, \theta)^{\frac{n}{\tau b}}. \quad (13)$$

In practice we absorb the exponent into a single constant $\beta := \frac{n}{\tau b}$ and tune β for performance. As per the hyperparameter details below, often $\beta < n$ is the most effective, indicating a *hot posterior* effect.

Monte Carlo estimates of $p(y_i \mid x_i, \theta)$ with surrogate models For a single LLM call $y = \mathbf{LLM}(x; \theta)$ that outputs an answer directly, computing $p(y \mid x, \theta)$ is as simple as summing log-probabilities across every token in y and summing them. However, complex LLM-based systems that include intermediate outputs and reasoning involve sources of stochasticity that are not captured in log-probabilities associated with the tokens of y .

Let z be a variable capturing all intermediate outputs in the computation of $y = \mathbf{LLM}(x; \theta)$. This includes internal LLM calls and reasoning text. We can express the generative process of sampling y as

$$z \sim p(z \mid x, \theta), \quad y \sim p(y \mid z, x, \theta). \quad (14)$$

Then the probability $p(y \mid x, \theta)$ required for MHLF would be computed as

$$p(y \mid x, \theta) = \sum_z p(y \mid z, x, \theta) p(z \mid x, \theta) = \mathbb{E}_{z \sim p(z \mid x, \theta)} [p(y \mid z, x, \theta)]. \quad (15)$$

As suggested by the second equality, this intractable integral is again amenable to Monte Carlo estimates. We use a single sample of z to estimate the likelihood in MHLF. We note that one alternative would be to remove all stochasticity from $z \sim p(z \mid x, \theta)$ by fixing a seed or setting the LLM temperature to 0 in the process of sampling z ; however, this would remove necessary stochasticity from the final model result and alter the underlying model.

When using closed-source model providers, log-probabilities and thus the value of $p(y \mid z, x, \theta)$ itself is sometimes withheld. In this case, during MHLF only we substitute the final LLM call in our LLM-based system with an open source model.

Lastly, we point out that all of the tricks above apply to computing probability masses for any LLM-based system, including the proposal density $q(\theta' \mid \theta)$ also required for MHLF.

B Experiment Details

B.1 Dataset Details

AIME [10], released under the MIT license, contains problems from the American Invitational Mathematics Examination (AIME)—a prestigious high school competition known for its challenging mathematical questions. Each answer is an integer. The exam consists of 29 to 30 questions per year. For evaluation, we used the 2024 exam, which was not included in GPT’s training data.

SimpleQA [63], released under the MIT license, is a benchmark that evaluates the ability of LLMs to answer short, fact-seeking questions. It covers a wide range of topics, including science, history, geography, history, politics, etc. Both its questions and answers are short and direct. In our experiments, we evaluated the models on a subset of 100 examples from the dataset.

QASPER [11], released under the CC-BY-4.0 license, is a free-form question-answering dataset focused on scientific research papers. It contains 5,049 questions across 1,585 papers in the field of Natural Language Processing. Each question is based on the content of a specific paper. In our experiments, we provided the model with a passage from the paper that contains the answer (i.e., the context), and then posed the question for it to answer using that context. We evaluated our model on 100 samples from this dataset under two different scenarios. In the first scenario, the context was entirely missing for 35 of the samples. In the second, 33 samples were provided with randomly selected context [65] that did not contain the correct answer. In both cases, the model was expected to abstain from answering.

B.2 Hyperparameters

In the following experiments, we use the OpenAI API for calls to GPT-4o-mini and GPT-4o. As our surrogate model for probability mass estimates (see Appendix A), we use Llama-3.1-Nemotron-70B-Instruct-HF [59, 3] through the Together AI API.

For all LLM calls we use a temperature of 1. We ensure $m = 10$ final answers are sampled for each method. For the Chain-of-Thought baseline, we used the initial prompt and sampled 10 answers, then aggregated the resulting answers. For TextGrad, we use the sample initial prompt but run TextGrad for a given number of steps before sampling 10 answers from the final prompt. For MCMC, we sample 10 individual prompts from a single MCMC chain and sample 1 answer from each. We tune the MHLF parameter β (see Appendix A) separately for each dataset. GPT-4o was employed for clustering and LLM-based evaluation. Further hyperparameter details are shown below and in our code (see especially the config files `qasper.yaml`, `aime.yaml`, and `simpleqa.yaml`).

Table 4: Hyperparameters used for each dataset and method

Dataset	Method	Model	Steps (T)	β	Burn-in (d)	Thinning (h)
AIME	Chain-of-Thought	GPT-4o	0	—	—	—
	TextGrad		60	—	—	—
	MHLF		60	10	6	6
SimpleQA	Chain-of-Thought	GPT-4o	0	—	—	—
	TextGrad		60	—	—	—
	MHLF		60	100	6	6
QASPER	Chain-of-Thought	GPT-4o-mini	0	—	—	—
	TextGrad		20	—	—	—
	MHLF		20	100	2	2

For all methods, we fix a string at the end of the prompt describing standardized formatting instructions for the model’s final answer. We extract this answer and evaluate likelihoods $p(y \mid z, x, \theta)$ only on this value, relegating any reasoning beforehand to the z variable (see Appendix A).

B.2.1 Baselines

We adapted the perturber baselines from SPUQ [17], specifically selecting the Paraphrasing and System Message perturbers for comparison. For all runs, we used GPT-4o-mini and GPT-4o for a fair comparison. Our implementation differs from the original in several details:

Paraphrasing: Rather than using a single LLM call with JSON formatting to produce all paraphrases, we made separate LLM calls for each paraphrase (to avoid invalid JSON outputs from the LLM with the original prompt). We used the following prompt:

```
Suggest a way to paraphrase the text in triple quotes above.
If the original text is a question, please make sure that your answer is also a question.
If the original text has answer options, please make sure your answer also has those options in the
    same order.
Answer should ONLY be the paraphrase and nothing else.
```

System Message: Instead of sampling with replacement from the available prompts, we expanded the set of system prompts and sampled without replacement. We appended these system prompts to the beginning of the message chain to preserve any existing system prompts. This was crucial for maintaining the output format required by the evaluator (e.g., answers ending with Answer: <THE ANSWER>). The set of system prompts used was:

```
"you are a helpful assistant"
"you are a question-answering assistant"
"you are a nice assistant"
"You are an AI support tool."
"You are a friendly helper."
"You are here to assist users."
"You provide useful answers."
"You are a kind AI agent."
"You offer good information."
"You are a smart assistant."
"You help with many tasks."
"You are a reliable AI."
"You give clear responses."
"You are an able assistant."
"You try to be useful."
"You are a positive AI."
"You guide users well."
"You are an adept helper."
"You simplify complex things."
"You are a virtual guide."
"You aim to be accurate."
```

B.2.2 Conformal Factuality

Factuality annotation Since there is no ground-truth output for the biography generation task, we assess the factuality of each generated answer by verifying its atomic sub-claims via web search. We use the GPT web search tool API⁶, which allows the model to retrieve external evidence before making a judgment. Each sub-claim is labeled as factual (1) or not (0) based on the retrieved information. We call the API as follows:

```
response = GPT_client.responses.create(
    model="gpt-4o-mini",
    tools=[
        {
            "type": "web_search_preview",
            "search_context_size": "low"
        }
    ],
    input=prompt,
)
response_content = response.output_text
```

Model and Prompt Setup We use GPT-4 for base biography generation and GPT-4o-mini for claim decomposition, factuality annotation, and frequency-based entailment scoring. Both the baseline frequency scoring and MHLP initialization use the same default system prompt: "You are a helpful assistant. Write a bio for people." For frequency scoring, we generate five alternative answers using this prompt. All prompts are listed in Table 5.

⁶<https://platform.openai.com/docs/guides/tools-web-search?api-mode=responses>

Subclaim Separator

Please breakdown the following input into a set of small, independent claims (make sure not to add any information), and return the output as a jsonl, where each line is subclaim:[CLAIM], gpt-score:[CONF].\n The confidence score [CONF] should represent your confidence in the claim, where a 1 is obvious facts and results like 'The earth is round' and '1+1=2'. A 0 is for claims that are very obscure or difficult for anyone to know, like the birthdays of non-notable people. If the input is short, it is fine to only return 1 claim. The input is:

Frequency scoring

You will get a list of claims and piece of text. For each claim, score whether the text supports, contradicts, or is unrelated to the claim. Directly return a jsonl, where each line is {"id":[CLAIM_ID], "score":[SCORE]}. Directly return the jsonl with no explanation or other formatting. For the [SCORE], return 1 for supports, -1 for contradicts, and 0 for unrelated. The claims are:\n{claim_string}\n\nThe text is:\n{output}

Factuality Annotation

Please verify if each of these claims is factual.\nClaims:\n[claims_text]\nReturn your answer as a JSON array, where each element is an object with these keys: {"subclaim": "[CLAIM]", "factual": 1 or 0, "source": "source or explanation"}\nFormat your response as a valid JSON array only, with no additional text or formatting.\nExample:\n[\n {\n "subclaim": "claim 1", "factual": 1, "source": "source"},\n {\n "subclaim": "claim 2", "factual": 0, "source": "source"}\n]\n

Table 5: Prompts for sub-claim separator, frequency scoring, and factuality annotation. Note both sub-claim separator and frequency scoring prompts are the same as used in [40]

Hyperparameter We run a single Metropolis-Hastings chain with $T = 20$ total steps, a burn-in of $d = 4$, and a thinning interval of $h = 4$, resulting in $m = 4$ sampled prompts. Together with the initial prompt, we obtain 5 prompts in total, which are used to compute frequency scores.

B.3 Examples

In this section, we exhibit several example questions and answers comparing results from Textual Bayes to TextGrad. We show how the the 10 answers sampled by each method are clustered and the number of answers that fall into each cluster. Overall, we see that Textual Bayes’s confidence levels are better calibrated to the model’s correctness.

B.3.1 SimpleQA

The following examples are selected from the SimpleQA dataset. The second example represents a case where the LLM appears truly to not know the answer; our method quantifies uncertainty better by expressing much lower confidence (40%) than the TextGrad baseline.

Question: According to Medland, Sarah E.; Loesch, Danuta Z.; Mdzewski, Bogdan; Zhu, Gu; Montgomery, Grant W.; Martin, Nicholas G. (September 28, 2007), what chromosome location was identified as linked to the finger ridge counts of the ring, index, and middle fingers through multivariate linkage analysis?

Answer: 5q14.1

Semantic Cluster	TextGrad	Ours
5q14.1	3	7
5q14.3	3	0
5	1	1
15q14	1	0
21q22	1	0
3q26	1	0
5q13	0	1
5q35	0	1

Question: What was the population of the town of Lesbury in Northumberland, England in the 2011 census?

Answer: 1007

Semantic Cluster	TextGrad	Ours
1,154	7	4
1,118	1	0
1,057	1	0
1,205	1	0
1,264	0	1
1,386	0	1
1,122	0	1
984	0	1
1,187	0	1
1,112	0	1

B.3.2 QASPER

The following examples are selected from the QASPER dataset.

Context: We begin with a hate speech lexicon containing words and phrases identified by internet users as hate speech, compiled by Hatebase.org. Using the Twitter API we searched for tweets containing terms from the lexicon, resulting in a sample of tweets from 33,458 Twitter users. We extracted the time-line for each user, resulting in a set of 85.4 million tweets. From this corpus we then took a random sample of 25k tweets containing terms from the lexicon and had them manually coded by CrowdFlower (CF) workers. Workers were asked to label each tweet as one of three categories: hate speech, offensive but not hate speech, or neither offensive nor hate speech. They were provided with our definition along with a paragraph explaining it in further detail. Users were asked to think not just about the words appearing in a given tweet but about the context in which they were used. They were instructed that the presence of a particular word, however offensive, did not necessarily indicate a tweet is hate speech. Each tweet was coded by three or more people. The intercoder-agreement score provided by CF is 92%. We use the majority decision for each tweet to assign a label. Some tweets were not assigned labels as there was no majority class. This results in a sample of 24,802 labeled tweets.

Question: How long is their dataset?

Answer: 85400000

Semantic Answer	TextGrad	Ours
85.4 million tweets	1	6
24,802 tweets	9	4

Random Context: Figure FIGREF4 is the overview of the proposed method using character 3-gram embeddings (char3-MS-vec). As illustrated in this figure, our proposed method regards the sum of char3-MS-vec and the standard word embedding as an input of an RNN. In other words, let $\text{char3-MS-vec} + \text{word-vec}$ and we replace Equation with the following: $\text{char3-MS-vec} + \text{word-vec}$

Question: Do they report results only on English data?

Answer: Unanswerable

Semantic Answer	TextGrad	Ours
Unclear / not specified in context	0	6
Results are only on English data	0	1
Results are not only on English	9	3
Formatting error in answer	1	0