

GENUINE: Graph Enhanced Multi-level Uncertainty Estimation for Large Language Models

Anonymous ACL submission

Abstract

Uncertainty estimation is essential for enhancing the reliability of Large Language Models (LLMs), particularly in high-stakes applications. Existing methods often overlook semantic dependencies, relying on token-level probability measures that fail to capture structural relationships within the generated text. We propose GENUINE: Graph ENhanced mUlti-level uncertaINty Estimation for Large Language Models, a structure-aware framework that leverages dependency parse trees and hierarchical graph pooling to refine uncertainty quantification. By incorporating supervised learning, GENUINE effectively models semantic and structural relationships, improving confidence assessments. Extensive experiments across NLP tasks show that GENUINE achieves up to 29% higher AUROC than semantic entropy-based approaches and reduces calibration errors by over 15%, demonstrating the effectiveness of graph-based uncertainty modeling. The code is available at <https://anonymous.4open.science/r/GUQ-39E7>.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in conversation (Wu et al., 2024), logical reasoning (Wang et al., 2023), and scientific discovery (Shojaee et al., 2024). Models such as GPT-4 (Achiam et al., 2023), Gemini (Team et al., 2023), and DeepSeek (Liu et al., 2024a), trained on vast corpora and aligned to human preferences, have significantly expanded the potential of AI. However, despite these advancements, LLMs are prone to well-documented reliability issues, including hallucinations and factual inaccuracies (Huang et al., 2025; Liu et al., 2024c). These issues pose serious risks, particularly in high-stakes applications such as medical diagnosis (Panagoulas et al., 2024), financial decision-making (de Zarzà et al., 2023), and legal advisory systems (Cheong et al., 2024), where

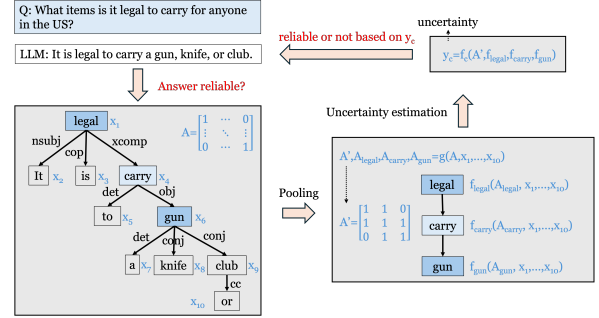


Figure 1: An example highlighting the role of graph pooling to identify tokens' significance in uncertainty estimation. Critical tokens are identified by graph pooling through dependency parsing tree and backpropagation of ground truth label, which makes the uncertainty task aware. A represents the adjacency matrix for a tree structure, where connected tree nodes are assigned value 1 while others are assigned value 0. g represents the pooling method, and f represents the information propagating through the pooling process.

users must rely on the model's outputs with confidence. Therefore, uncertainty quantification(UQ), which assesses the trustworthiness of an LLM response, is essential for safe and effective human and artificial intelligence interaction.

UQ in LLM-generated outputs presents several challenges. First, LLMs often produce long-form textual responses, making attributing uncertainty to specific components difficult. Second, uncertainties may affect only a few critical tokens within an otherwise coherent response, undermining the reliability of the entire output. Third, aggregating uncertainty across multiple tokens in lengthy outputs is non-trivial, requiring distinguishing semantically pivotal tokens from those not pivotal.

Previous studies have explored various approaches to quantify uncertainty in LLM outputs. Some methods rely on self-evaluation through modified prompts (Tian et al., 2023b), though they often inherit the model's biases. Others use token-level uncertainty measures based on logits, entropy, or probability distributions (Kuhn et al., 2023; Ma-

linin and Gales, 2020, 2021). Recent advancements, such as semantic entropy, cluster semantically equivalent generations and measure entropy as an uncertainty indicator (Kuhn et al., 2023). However, most existing methods treat all tokens equally, overlooking findings that certain tokens carry more semantic weight in determining output validity (Liu et al., 2024b; Duan et al., 2024; Cheng and Vlachos, 2024). Additionally, some approaches (Duan et al., 2024) depend on external smaller models to estimate token importance, but these models often operate independently of the LLM’s internal representations. As a result, they may introduce inconsistencies, misinterpret token dependencies, or fail to capture the structural relationships within the generated text, leading to inaccurate uncertainty estimates.

To illustrate this issue, consider the example in Fig. 1. A user inquires about legal items to carry in the United States, but the model responds with a list of illegal items, such as a gun, knife, and club. The misunderstanding stems from the token "legal," which is central to the query’s meaning. A minor modification, replacing the word legal with illegal, would render the response appropriate. This example underscores two insights: certain tokens are disproportionately influential in determining output validity. Dependency parse trees effectively capture the hierarchical structure of sentence meaning by identifying core decision points. Building on these insights, we propose leveraging dependency parse trees and graph pooling techniques to infer LLM prediction uncertainty in a structured and interpretable manner.

Modeling uncertainty estimation as a graph-based problem offers several advantages. Graphs inherently capture dependencies between generated tokens, reflecting the autoregressive nature of LLMs, where each token influences subsequent ones. By representing an LLM response as a structured graph, we can propagate and aggregate critical information across tokens, ensuring that semantically significant tokens contribute more substantially to the overall uncertainty estimate. However, this approach introduces several challenges. Determining the optimal graph structure that accurately represents token dependencies remains an open question. Selecting appropriate graph pooling techniques that summarize uncertainty information effectively without losing essential context is difficult. Addressing these challenges is essential to fully realize the potential of graph-based uncer-

tainty estimation.

Our approach integrates multiple uncertainty features to enhance robustness. Specifically, we utilize probability distributions, entropy-based measures, and LLM embeddings to model uncertainty. We introduce a hierarchical strategy to address the challenge of aggregating uncertainty over long-form text. We construct a dependency parse tree for each sentence to extract structural and semantic relationships. We merge sentence-level trees into a document-level graph by connecting their root nodes. We apply graph pooling techniques to model uncertainty across the entire paragraph efficiently. GENUINE involves learning pooling functions that adaptively fuse different features, capturing both local and global dependencies within the text. Experimental results prove that GENUINE outperforms other baselines, highlighting the critical role of structural relationships in uncertainty estimation. Furthermore, we compare the effectiveness of probability-based and embedding-based features across various datasets and LLMs, offering insights into their respective utilities. Given that commercial LLMs typically provide only probability-related features, our findings suggest an intriguing direction for future research. Exploring whether open-source LLMs, which offer both probability and embedding features, can facilitate superior UQ compared to their commercial counterparts.

The following are our main contributions:

- We highlight the role of semantically significant tokens in uncertainty estimation, demonstrating how structural relationships can enhance model uncertainty assessment.
- We propose a graph-based framework for LLM UQ, integrating dependency parse trees and graph pooling to capture structural and semantic relationships in the generated text.
- We develop an adaptive graph pooling mechanism that effectively propagates and aggregates uncertainty information by learning to fuse multiple uncertainty features.
- We conduct extensive experiments on real-world datasets, demonstrating that GENUINE outperforms existing UQ methods in assessing the trustworthiness of LLM-generated responses.

2 Related Works

Uncertainty Quantification in LLMs. Uncertainty quantification is well-studied in traditional machine learning (Chen et al., 2019; Zhao et al.,

2020), but remains challenging for LLMs due to their open-ended outputs, where multiple valid responses can exist. This flexibility complicates uncertainty estimation, requiring methods beyond standard predictive confidence. Current approaches fall into two categories. Self-assessment prompts LLMs to estimate their own uncertainty (Kadavath et al., 2022; Lin et al., 2022; Tian et al., 2023a), but often reflect model biases and inconsistencies. External methods assess uncertainty via output consistency (Manakul et al., 2023) or entropy measures (Malinin and Gales, 2020), though these typically assume uniform token importance, overlooking the fact that certain tokens contribute more to the overall reliability of a response. Recent work addresses this limitation by incorporating semantic awareness. Semantic entropy (SE) (Kuhn et al., 2023) reduces redundancy by grouping semantically equivalent outputs. Others re-weight token contributions (Duan et al., 2024) or leverage hidden activations as uncertainty signals (Liu et al., 2024b). Building on this, we integrate dependency parse trees to identify key tokens shaping response meaning, while hidden activations provide semantic context. This combination enables a structured and context-aware approach to uncertainty estimation in LLMs.

Graph Pooling Approaches. Graph pooling condenses input graphs while preserving key structural and semantic information. It generally falls into flat pooling, which applies simple aggregation functions like mean or sum (Xu et al., 2019; Duvenaud et al., 2015), and hierarchical pooling, which progressively coarsens the graph to capture multi-level relationships (Ying et al., 2018). Notable hierarchical methods include DiffPool (Ying et al., 2018), which learns adaptive pooling assignments, and StructPool (Yuan and Ji, 2020), which incorporates high-order structural dependencies. Other strategies include memory-based pooling (Khasahmadi et al., 2020), spectral filtering (Defferrard et al., 2016), and expressive pooling architectures (Bianchi and Lachi, 2023). Unsupervised pooling techniques, like mutual information maximization (Liu et al., 2022), further enable structure-preserving and label-free compression. This work proposes a hierarchical pooling approach leveraging dependency tree structures to improve uncertainty estimation. By representing LLM outputs as dependency graphs, GENUINE captures both semantic and structural relationships, prioritizing key tokens for a more accurate and

interpretable uncertainty assessment.

3 Background

This section defines the problem, provides the necessary background, and features helpful for uncertainty estimation in LLMs, laying the foundation for our proposed approach.

3.1 Problem Setup

Uncertainty quantification in LLMs involves assessing confidence in LLM-generated responses based on input prompts. Given a prompt $\mathbf{x} = \{x_1, x_2, \dots, x_k\}$, an LLM generates an output sequence $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$, where each token y_j is sampled from a probability distribution conditioned on the prompt and prior tokens:

$$y_j \sim p_\theta(\cdot | \mathbf{x}, y_1, y_2, \dots, y_{j-1}), \quad (1)$$

where p_θ represents the model’s learned parameters. This next-token probability reflects how likely the model is to generate a particular token given the preceding context. Following (Liu et al., 2024b), when there is a downstream task, such as question answering or machine translation, a scoring function is introduced to evaluate the quality of the generated output. For such kinds of evaluation functions, factual truth or humans usually decide the true response. Thus the uncertainty estimation task can be framed as a function $g(\mathbf{x}, \mathbf{y})$ that predicts the expected correctness of a response:

$$g(\mathbf{x}, \mathbf{y}) \approx \mathbb{E}[s(\mathbf{y}, \mathbf{y}_{\text{true}}) | \mathbf{x}, \mathbf{y}]. \quad (2)$$

Here, $s(\mathbf{y}, \mathbf{y}_{\text{true}})$ denotes an evaluation metric comparing the generated response \mathbf{y} with a ground-truth reference \mathbf{y}_{true} . The expectation is taken considering the semantic flexibility of natural language. The uncertainty arises from the input prompt \mathbf{x} and the LLM itself rather than from a single absolute reference answer.

3.2 Dependency Parse Trees in NLP

Dependency parse trees provide a structured representation of syntactic relationships, defining hierarchical dependencies such as *subjects*, *objects*, and *modifiers* within a sentence. These structures have been widely applied in various NLP tasks, including relation extraction (RE) (Fundel et al., 2006; Björne et al., 2009), named entity recognition (NER) (Jie et al., 2017), and semantic role labeling (SRL) (Marcheggiani and Titov, 2017). They also enhance summarization by prioritizing

salient information while filtering redundant content (Li et al., 2014; Xu and Durrett, 2019). This work uses dependency parse trees to model structural relationships in LLM-generated text. These trees serve two key purposes: (1) They provide a hierarchical organization of tokens, helping distinguish pivotal words that shape response meaning, (2) They offer a consistent structure across different sentence formations, making them adaptable for modeling uncertainty in diverse LLM outputs.

3.3 Features for Uncertainty Estimation

Uncertainty estimation in LLMs relies on extracting meaningful features from the generated text. Prior studies (Xiao et al., 2022; Kadavath et al., 2022; Lin et al., 2022; Tian et al., 2023a; Kuhn et al., 2023; Liu et al., 2024b) have demonstrated the effectiveness of token-level probability metrics. We categorize these features based on their sources (Liu et al., 2024b):

White-box features: These features are derived from hidden-layer activations, capturing the internal representation of tokens and providing insights into model confidence. These features are available only in open-source LLMs.

Grey-box features: These include *token probabilities* and transformations such as entropy, offering uncertainty signals applicable to both open-source and commercial LLMs. The entropy of a discrete distribution p over the vocabulary \mathcal{V} is defined as $H(p) = -\sum_{v \in \mathcal{V}} p(v) \log(p(v))$. Given a prompt-response pair $(\mathbf{x}, \mathbf{y}) = (x_1, \dots, x_k, y_1, \dots, y_n)$, the entropy features for the j -th output token are given by $H(q_\theta(y_j | \mathbf{x}, y_1, \dots, y_{j-1}))$, where q_θ denotes the LLM. The detailed mathematical definition of the features is provided in Appendix A.2.

4 Approach

This section details our approach, including graph formulation, hierarchical learning, and joint optimization, enabling a more structured and context-aware uncertainty estimation for LLMs.

4.1 Graph Formulation

We transform dependency parse trees into graphs to structure LLM-generated text for uncertainty estimation. We first obtain the dependency tree using the Stanford NLTK parser, where each word serves as a node, and directed edges represent dependency relations. As shown in Fig. 2, the root word, such as "prefer," has dependent words like "I" and "flight," forming a tree-like structure.

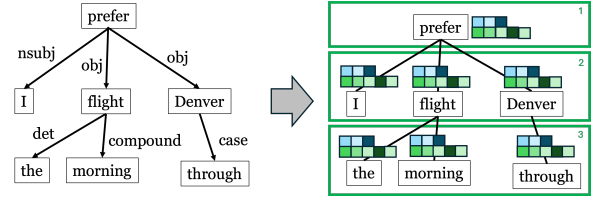


Figure 2: Dependency parse tree example. Each tree node is one token from the output. If two tokens have a relation, they are connected. Each tree node has additional features, such as probability, entropy, and embeddings(only for open-box LLMs).

To extend this formulation beyond individual sentences, we construct a paragraph-level graph by linking the root nodes of multiple sentence-level dependency trees. Prior work (Duan et al., 2024) estimates uncertainty at the sentence level using a separate model to compute similarity, but such approaches may overlook deeper semantic relationships between sentences. Instead, GENUINE learns inter-sentence relations directly, ensuring a more cohesive uncertainty estimation. Connecting root nodes across sentences enables cross-sentence token interactions, allowing uncertainty information to propagate effectively across the entire output. This formulation ensures that pivotal words influence the overall confidence estimation. The resulting global dependency graph provides a structured representation of LLM output, enhancing the ability of the proposed approach to assess uncertainty in LLM-generated text.

4.2 Hierarchical Learning

Transforming dependency parse trees into graphs enables us to frame uncertainty estimation as a graph aggregation problem, where each LLM-generated output is represented as a graph with nodes corresponding to words and edges capturing dependency relations. Each node has token-level features, such as next-token probability, entropy, and hidden state embeddings. We propose a hierarchical graph pooling approach inspired by semantic parsing trees (Song and King, 2022) to aggregate this information efficiently.

In a dependency graph (Fig. 2), words appear at different levels based on their distance from the root token, which often signifies their semantic importance. Higher-level words generally play a more critical role in defining the sentence’s meaning and, consequently, have a greater impact on uncertainty. To capture this, we introduce graph pooling, which groups tokens at different hierarchical levels, mitigating the effect of noisy words while assigning appropriate contributions to each

token’s uncertainty estimate.

Formally, given a dependency graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} represents words and \mathcal{E} defines their syntactic relations, we define an adjacency matrix $\mathcal{A} \in \mathbb{R}^{n \times n}$ and a feature matrix $\mathcal{X} \in \mathbb{R}^{n \times d}$. Inspired by hierarchical graph pooling methods (Ying et al., 2018), we define the node clustering process using a learned soft assignment matrix:

$$\mathcal{S}^l = \text{Softmax}(f(\mathcal{A}^l, \mathcal{X}^l, \theta_s)), \quad (3)$$

where \mathcal{A}^l and \mathcal{X}^l represent the adjacency and feature matrices at pooling layer l , and f in a GNN with learnable parameters θ_s .

Before pooling, information propagates across the graph to model connectivity between clusters:

$$\mathcal{Z}^l = f(\mathcal{A}^l, \mathcal{X}^l, \theta_z), \quad (4)$$

where θ_z are the parameters of the GNN responsible for feature transformation. Using the learned assignment matrix \mathcal{S}^l , the graph is iteratively coarsened to generate a more compact representation:

$$\begin{aligned} \mathcal{X}^{l+1} &= \mathcal{S}^l \mathcal{Z}^l \in \mathbb{R}^{n_{l+1} \times d}, \\ \mathcal{A}^{l+1} &= \mathcal{S}^l \mathcal{A}^l \mathcal{S}^{lT} \in \mathbb{R}^{n_{l+1} \times n_{l+1}}. \end{aligned} \quad (5)$$

Here, \mathcal{X}^l and \mathcal{A}^l are iteratively refined representations at each pooling level, ensuring that semantically important tokens retain greater influence. By hierarchically aggregating token-level uncertainty, GENUINE enhances interpretability and robustness, providing a structured estimation of confidence in LLM-generated responses.

4.3 Joint Optimization

Uncertainty estimation in LLMs relies on multiple features, as discussed in Section 3.3, including hidden states (white-box features) and probability-based signals (grey-box features), each contributing differently. Prior work (Liu et al., 2024b) shows that hidden states encode valuable uncertainty information, partly due to the misalignment between pretraining objectives and uncertainty estimation. Moreover, hidden states capture semantic relationships among tokens, making them especially important for confidence evaluation.

We propose a joint optimization framework to effectively integrate multiple uncertainty features. As illustrated in Fig. 3, GENUINE includes a semantic pooling module that leverages hidden state embeddings and a structural pooling module that utilizes probability and entropy features. Both modules operate on a shared dependency parse tree,

providing a unified structural backbone. Their outputs are combined via a fusion module that learns a joint graph pooling matrix, balancing semantic and structural signals to refine uncertainty estimation. Instead of merging features at the node level, we fuse them at the assignment matrix level to better balance structural and semantic information. This design is motivated by three factors. First, direct feature fusion would bias toward embeddings due to their higher dimensionality. Second, embeddings encode semantic context but lack precise generation uncertainty, while probability and entropy features provide more accurate confidence signals. Third, the assignment matrix inherently reflects token importance and relational structure, making it a more effective fusion point for heterogeneous features.

To achieve this, we introduce an end-to-end learnable fusion module, where the fused assignment matrix is computed as:

$$\mathcal{S}_*^l = \text{Softmax}(g(\mathcal{S}_{\text{grey}}^l, \mathcal{S}_{\text{white}}^l, \theta_{s*})), \quad (6)$$

where $\mathcal{S}_{\text{grey}}^l$ and $\mathcal{S}_{\text{white}}^l$ are the assignment matrices at pooling layer l from the structural and semantic modules, respectively, and θ_{s*} denotes the learnable parameters of the fusion function g .

Following this, a GNN propagates information across the graph, refining node representations:

$$\mathcal{Z}_*^l = f(\mathcal{A}_*^l, \mathcal{X}_*^l, \theta_{z*}), \quad (7)$$

where f is a GNN with learnable parameters θ_{z*} . These updated assignment and node embedding matrices are used to refine the graph iteratively:

$$\begin{aligned} \mathcal{X}_*^{l+1} &= \mathcal{S}_*^l \mathcal{Z}_*^l, \\ \mathcal{A}_*^{l+1} &= \mathcal{S}_*^l \mathcal{A}_*^l \mathcal{S}_*^{lT}. \end{aligned} \quad (8)$$

Here, \mathcal{X}_* encodes probability and entropy features, while embeddings enhance the model’s semantic understanding. The independent assignment matrices $\mathcal{S}_{\text{grey}}^l$ and $\mathcal{S}_{\text{white}}^l$ are jointly optimized to capture both structural and contextual uncertainty, improving the robustness of LLM confidence evaluation. Recall that when using open-box LLMs, which allow users access to grey-box features and white-box features, our fusion process can be directly applied. While using black-box LLMs, which only allow access to grey-box features, the fusion process can not proceed without white-box features. However, this will not hinder the application of GENUINE as the graph structures and the joint

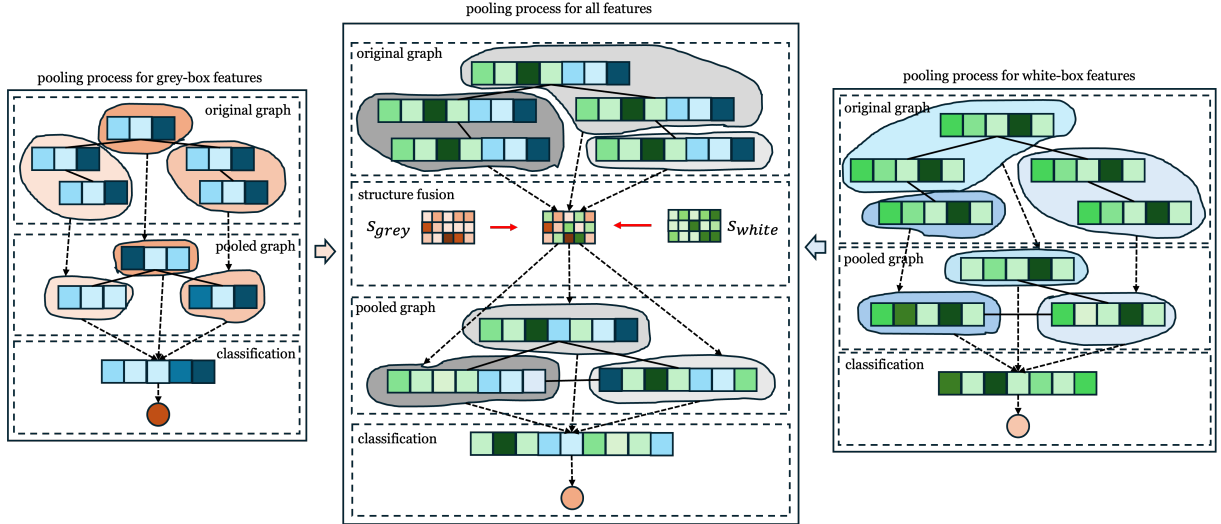


Figure 3: The Overview of GENUINE, composed of three modules: (1) pooling based on grey-box features, (2) pooling based on white-box features, and (3) a learnable fusion process integrating both modules. Both the grey-box pooling process and white-box pooling process share the same graph structure, but differ in features, which leads to different fusion matrices. The structure fusion process helps better integrate various fusion matrices.

optimization will remain effective in estimating the uncertainty. We provide both the results with only grey-box features and both grey-box and white-box features in the experiments in Fig. 4 to prove.

5 Experiments

This section evaluates GENUINE across multiple dimensions: (1) effectiveness in assessing uncertainty(Section 5.2), (2) an ablation study to analyze the role of two modules(Section 5.3), (3) a scalability test to assess computational efficiency(Section 5.4), (4) the impact of dependency parse trees on uncertainty estimation (Appendix B.2), (5) a parameter analysis to determine the sensitivity of GENUINE to hyperparameter tuning(Appendix B.4), (6) the impact of LLM parameters on GENUINE’s uncertainty estimation performance(Appendix B.5), and (7) the impact of training dataset size and noisy labels on GENUINE’s performance(Section 5.5 and Appendix B.6). Due to space constraints, the results on dimensions 4, 5, 6, and 7 are presented in Appendix B.

5.1 Experimental Setup

We evaluate GENUINE using different LLM architectures, multiple datasets spanning various NLP tasks, and state-of-the-art baselines.

LLMs. We consider open-source LLMs, including Llama2-7B, Llama2-13B, Llama3-8B (Touvron et al., 2023), as well as Gemma-7B and Gemma2-9B (Gemma Team et al., 2024). The respective tokenizers provided by Hugging Face are used, and model parameters remain unchanged.

Datasets. We evaluate uncertainty estimation on three NLP tasks: question answering(CoQA(Reddy et al., 2019), TriviaQA (Joshi et al., 2017), and Finance QA dataset (Taori et al., 2023)), machine translation(WMT 2014 dataset (Bojar et al., 2014)), and summarization(CNN dataset (Hermann et al., 2015)). The details of the datasets is introduced in Appendix A.1. Each dataset is split into training (60%), validation (10%), and test (30%) sets, with five runs performed to mitigate the effects of randomness in parameter optimization. Few-shot prompting is adopted, with templates detailed in Appendix A.3.

Baselines. We include five categories of state-of-the-art baselines to compare GENUINE against with: (1) A4C (Tian et al., 2023b), which directly queries the LLM for its self-assessed uncertainty, (2) Entropy and probability-based methods, including Avg Probability (Prob) and Avg Entropy (Ent), as defined in Table 3 in the Appendix A.2, (3) Semantic-aware methods, such as Semantic Entropy (SE) (Kuhn et al., 2023) and SAR (Duan et al., 2024), (4) Bayesian based methods, including BayesPE(Tonolini et al., 2024), and (5) A supervised uncertainty estimation approach (Sup) (Liu et al., 2024b). Details of the prompt templates are provided in the Appendix A.3.

Evaluation Metrics. Following (Liu et al., 2024b; Kuhn et al., 2023), we evaluate GENUINE’s ability to distinguish correct from incorrect responses using uncertainty scores. Our primary metric is AU-ROC, which measures how well the model ranks correct responses above incorrect ones. We also

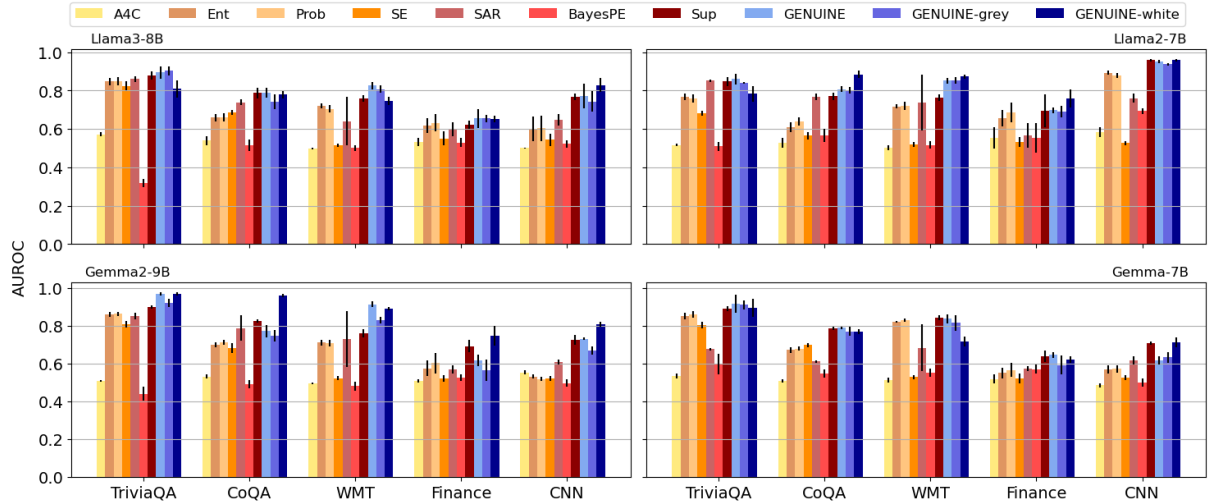


Figure 4: Comparison of AUROC on five datasets, four LLMs, and seven baselines. Error bars denote variance over five runs. GENUINE and its transformations outperform other baselines for all five datasets and four LLMs. Especially for datasets with relatively longer output from LLMs, such as WMT, Finance, and CNN datasets

assess calibration using Expected Calibration Error(ECE)(Naeini et al., 2015), and report Brier score(Hernández-Orallo et al., 2011) and negative log-likelihood(NLL)(Hastie et al., 2001) to evaluate the alignment between predicted uncertainty and true confidence. AUROC are shown in the main paper, others are reported in Appendix B.1.

5.2 Performance of Uncertainty Estimation

We evaluate GENUINE using the AUROC metric against state-of-the-art baselines. As shown in Fig. 4, GENUINE consistently outperforms prior methods, particularly on long-form generation tasks (WMT, Finance, CNN). Its dependency-based structural modeling improves uncertainty estimation by reducing error propagation across extended sequences. GENUINE also achieves better calibration, as evidenced by lower ECE, NLL, and Brier scores (Appendix B.1), minimizing uncertainty misalignment in downstream tasks. The results further highlight that response length significantly impacts uncertainty estimation. As detailed in Table 5, GENUINE offers modest AUROC gains on shorter outputs (e.g., TriviaQA, CoQA), but shows substantial improvements on longer responses (e.g., WMT, Finance, CNN). Traditional token-wise methods accumulate errors over extended text, whereas GENUINE’s structured approach better handles long-form content, critical for tasks like dialogue and summarization.

Feature selection also plays a crucial role in uncertainty estimation. While combining multiple features generally improves performance, hidden-layer embeddings alone (GENUINE-white) perform best on Finance and CNN datasets, where

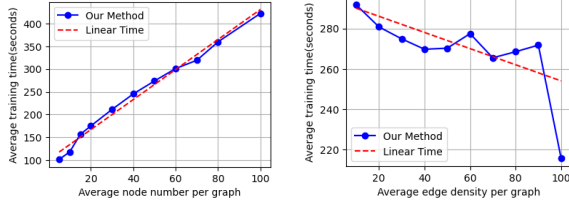
longer sequences amplify token-level error in entropy-based methods. To support both black-box and open-box LLMs, we introduce two variants: GENUINE-grey (using only grey-box features) and GENUINE-white (white-box features). The results demonstrate that in most cases, GENUINE-grey still has the superiority of performance, which shows the applicability of GENUINE in black-box LLMs. These findings also highlight the advantage of open-source LLMs with access to internal representations for robust uncertainty modeling.

5.3 Ablation Study

GENUINE introduces a graph structure and fusion mechanism to balance structural and semantic information. We conduct ablation studies on the TriviaQA dataset using Llama3-8B, Llama2-7B, Gemma2-9B, and Gemma-7B to assess the contribution of each component. Due to space constraints, we report results for Llama3-8B and Gemma2-9B, with full results in Appendix B.3. We denote variants as GENUINE w/o fusion & graph (without both modules) and GENUINE w/o fusion (with graph, but no fusion). As shown in Table 1, the graph structure and the fusion process improve AUROC on all the LLMs we use in our experiments. These findings highlight the graph structure and the fusion strategy’s effectiveness in integrating structural and semantic signals, enabling better uncertainty propagation. We observe that the improvement in Gemma2-9B model is not significant. The smaller gains for Gemma2-9B may be due to its already strong baseline performance (AUROC > 0.95), leaving limited room for improvement.

Table 1: Ablation study on TriviaQA. GENUINE w/o fusion & graph (without both modules) and GENUINE w/o fusion (with graph, but no fusion) (\uparrow means the higher the better)

Methods	Llama3-8B	Gemma2-9B
	AUROC \uparrow	AUROC \uparrow
GENUINE w/o fusion & graph	0.789 \pm 0.031	0.956 \pm 0.002
GENUINE w/o fusion	0.809 \pm 0.096	0.963 \pm 0.015
GENUINE	0.894 \pm 0.032	0.969 \pm 0.009



(a) Scalability test on the number of nodes per graph. (b) Scalability test on graph density.

Figure 5: Scalability test on the node number and edge density

5.4 Scalability

We evaluate the scalability of GENUINE by examining its computational efficiency with increasing node count and graph density. As shown in Fig. 5a, training time scales near-linearly with the number of nodes, demonstrating that GENUINE remains computationally feasible even for larger graphs. This suggests that the model can efficiently process uncertainty in large-scale LLM outputs without excessive overhead. In Fig. 5b, computational cost decreases as graph density increases, indicating that denser graphs facilitate more efficient uncertainty aggregation. Sparse graphs (e.g., 10% density) require 1.5 times more processing time than fully connected graphs (100% density), emphasizing the trade-off between structure complexity and efficiency. These findings confirm that GENUINE scales effectively with increasing graph complexity, making it well-suited for high-dimensional NLP tasks such as document summarization, multi-turn dialogue, and knowledge-intensive reasoning. Its ability to maintain efficiency while capturing semantic and structural relationships ensures its adaptability to real-world LLM evaluation scenarios.

5.5 Robustness

In real-world scenarios, uncertainty estimation models often face limited training data and noisy labels, which can affect performance. To evaluate the robustness of GENUINE under such conditions, we conduct experiments using the Llama3-8B model on the TriviaQA dataset. Table 8 in the

Appendix B.6 shows how varying training set sizes impact performance. Please refer to the Appendix for more details. While Table 2 examines the effect of label noise. For the latter, we randomly corrupt a portion of training labels (as specified by the noise ratio) and assess performance on the clean test set. Specifically, Table 2 shows that label noise negatively affects model performance. But GENUINE remains robust when up to 0.1% of the training labels are corrupted. However, AUROC declines sharply when the noise ratio increases, and so do the calibration metrics. These experiments demonstrate GENUINE’s resilience to data scarcity and label noise, highlighting its applicability in real-world settings.

Table 2: The impact of noisy labels on GENUINE performance. GENUINE remains robust with 0.1% of labels being noisy. (\uparrow means the higher the better, \downarrow means the lower the better)

noise ratio	AUROC \uparrow	ECE \downarrow
0	0.894 \pm 0.032	0.246 \pm 0.007
0.001	0.894 \pm 0.017	0.244 \pm 0.010
0.003	0.863 \pm 0.033	0.243 \pm 0.009
0.005	0.855 \pm 0.024	0.243 \pm 0.013
0.01	0.821 \pm 0.014	0.240 \pm 0.014
0.02	0.705 \pm 0.037	0.235 \pm 0.015
0.03	0.746 \pm 0.095	0.232 \pm 0.019
0.04	0.672 \pm 0.140	0.234 \pm 0.022
noise ratio	NLL \downarrow	Brier \downarrow
0	0.362 \pm 0.005	0.094 \pm 0.002
0.001	0.364 \pm 0.005	0.095 \pm 0.002
0.003	0.366 \pm 0.002	0.096 \pm 0.001
0.005	0.370 \pm 0.008	0.098 \pm 0.004
0.01	0.377 \pm 0.009	0.101 \pm 0.004
0.02	0.390 \pm 0.019	0.107 \pm 0.009
0.03	0.407 \pm 0.026	0.114 \pm 0.012
0.04	0.414 \pm 0.031	0.117 \pm 0.014

6 Conclusion

This paper introduces dependency-based semantic structures for uncertainty estimation in LLMs. Our findings prove that incorporating structural information enhances uncertainty modeling, leading to more accurate and calibrated estimates. GENUINE outperforms existing uncertainty estimation methods (AUROC), particularly in long-form text generation, while also improving calibration metrics (ECE, NLL, Brier). Our results show that semantic graphs derived from dependency parse trees enhance uncertainty modeling, making them valuable for evaluating LLMs’ outputs and guiding future improvements in adaptive uncertainty estimation in dynamic, real-world settings.

7 Ethical Consideration

GENUINE enhances the credibility and reliability of LLMs by improving uncertainty estimation, helping to mitigate the risks of misinformation. By refining confidence assessment, GENUINE reduces misinformation and promotes more trustworthy AI-generated content.

However, several ethical limitations must be considered. Uncertainty estimation does not prevent misinformation but provides a measure of confidence, which still requires human interpretation. Over-reliance on uncertainty scores could lead to misjudgments, either overestimating or underestimating the reliability of LLM outputs. Additionally, GENUINE’s effectiveness depends on dependency parsing and feature selection, which may introduce biases if trained on imbalanced datasets. Furthermore, while GENUINE improves model calibration, uncertainty quantification remains imperfect, and its reliability may vary across domains, particularly in high-stakes applications such as healthcare, finance, and law. Addressing these challenges requires ongoing evaluation, transparency, and responsible deployment to ensure ethical and fair AI use.

8 Limitations

GENUINE introduces a graph-based approach for confidence evaluation in LLMs, but certain limitations remain. GENUINE relies on token logits and embeddings, which, though widely available in open-source and commercial LLMs, may limit its applicability in black-box scenarios where such information is restricted. Additionally, its performance is influenced by generation length and labeled data availability, making it sensitive to dataset variability. Finally, this study focuses on NLP tasks and datasets, leaving open the exploration of its effectiveness in multimodal and cross-domain applications.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Filippo Maria Bianchi and Veronica Lachi. 2023. The expressive power of pooling in graph neural networks. *Advances in neural information processing systems*, 36:71603–71618.

Jari Björne, Juho Heimonen, Filip Ginter, Antti Airola, Tapio Pahikkala, and Tapio Salakoski. 2009. [Extracting complex biological events with rich graph-based feature sets](#). In *Proceedings of the BioNLP 2009 Workshop Companion Volume for Shared Task*.

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

Xuelu Chen, Muhao Chen, Weijia Shi, Yizhou Sun, and Carlo Zaniolo. 2019. Embedding uncertain knowledge graphs. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3363–3370.

Julius Cheng and Andreas Vlachos. 2024. [Measuring uncertainty in neural machine translation with similarity-sensitive entropy](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2115–2128, St. Julian’s, Malta. Association for Computational Linguistics.

Inyoung Cheong, King Xia, KJ Kevin Feng, Quan Ze Chen, and Amy X Zhang. 2024. (a) i am not a lawyer, but...: Engaging legal experts towards responsible llm policies for legal advice. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 2454–2469.

I de Zarzà, J de Curtò, Gemma Roig, and Carlos T Calafate. 2023. Optimized financial planning: integrating individual and cooperative budgeting models with llm recommendations. *AI*, 5(1):91–114.

Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29.

Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2024. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5050–5063.

David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gómez-Bombarelli, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *Proceedings of the 29th International Conference on Neural Information Processing Systems - Volume 2, NIPS’15*, page 2224–2232, Cambridge, MA, USA. MIT Press.

Katrin Fundel, Robert Küffner, and Ralf Zimmer. 2006. [RelEx – Relation extraction using dependency parse trees](#). *Bioinformatics*.

745	Thomas Mesnard Gemma Team, Cassidy Hardin,	Chin-Yew Lin and Franz Josef Och. 2004. Auto-	801
746	Robert Dadashi, Surya Bhupatiraju, Laurent Sifre,	matic evaluation of machine translation quality using	802
747	Morgane Rivi�re, Mihir Sanjay Kale, Juliette Love,	longest common subsequence and skip-bigram statis-	803
748	Pouya Tafti, L�onard Hussenot, and et al. 2024.	tics . In <i>Proceedings of the 42nd Annual Meeting of</i>	804
749	Gemma .	<i>the Association for Computational Linguistics (ACL-</i>	805
		<i>04)</i> , pages 605–612, Barcelona, Spain.	806
750	Trevor Hastie, Robert Tibshirani, and Jerome Friedman.	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	807
751	2001. <i>The Elements of Statistical Learning</i> . Springer	Teaching models to express their uncertainty in	808
752	Series in Statistics. Springer New York Inc., New	words. <i>Transactions on Machine Learning Research</i> .	809
753	York, NY, USA.		
754	Karl Moritz Hermann, Tomas Kocisky, Edward Grefen-	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	810
755	stette, Lasse Espeholt, Will Kay, Mustafa Suleyman,	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	811
756	and Phil Blunsom. 2015. Teaching machines to read	Deng, Chenyu Zhang, Chong Ruan, et al. 2024a.	812
757	and comprehend. <i>Advances in neural information</i>	Deepseek-v3 technical report. <i>arXiv preprint</i>	813
758	<i>processing systems</i> , 28.	<i>arXiv:2412.19437</i> .	814
759	Jos� Hern�ndez-Orallo, Peter A Flach, and C�sar Ferri	Linyu Liu, Yu Pan, Xiaocheng Li, and Guanting Chen.	815
760	Ramirez. 2011. Brier curves: a new cost-based vi-	2024b. Uncertainty estimation and quantification for	816
761	visualisation of classifier performance. In <i>Icml</i> , pages	llms: A simple supervised approach. <i>arXiv preprint</i>	817
762	585–592.	<i>arXiv:2404.15993</i> .	818
763	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	Ning Liu, Songlei Jian, Dongsheng Li, and Hongzuo Xu.	819
764	Zhangyin Feng, Haotian Wang, Qianglong Chen,	2022. Unsupervised hierarchical graph pooling via	820
765	Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting	substructure-sensitive mutual information maximiza-	821
766	Liu. 2025. A survey on hallucination in large lan-	tion . In <i>Proceedings of the 31st ACM International</i>	822
767	guage models: Principles, taxonomy, challenges, and	<i>Conference on Information & Knowledge Manage-</i>	823
768	open questions . <i>ACM Trans. Inf. Syst.</i> , 43(2).	<i>ment, CIKM ’22</i> , page 1299–1308, New York, NY,	824
769	Zhanming Jie, Aldrian Obaja Muis, and Wei Lu. 2017.	USA. Association for Computing Machinery.	825
770	Efficient dependency-guided named entity recogni-	Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying	826
771	tion. In <i>Thirty-First AAAI Conference on Artificial</i>	Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov,	827
772	<i>Intelligence</i> .	Muhammad Faaiz Taufiq, and Hang Li. 2024c. Trust-	828
773	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke	worthy llms: a survey and guideline for evaluat-	829
774	Zettlemoyer. 2017. TriviaQA: A large scale distant-	ing large language models’ alignment . <i>Preprint</i> ,	830
775	supervised challenge dataset for reading comprehen-	<i>arXiv:2308.05374</i> .	831
776	sion . In <i>Proceedings of the 55th Annual Meeting of</i>	Andrey Malinin and Mark Gales. 2020. Uncertainty	832
777	<i>the Association for Computational Linguistics (Vol-</i>	estimation in autoregressive structured prediction. In	833
778	<i>ume 1: Long Papers)</i> , pages 1601–1611, Vancouver,	<i>International Conference on Learning Representa-</i>	834
779	Canada. Association for Computational Linguistics.	<i>tions</i> .	835
780	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	Andrey Malinin and Mark Gales. 2021. Uncertainty	836
781	Henighan, Dawn Drain, Ethan Perez, Nicholas	estimation in autoregressive structured prediction. In	837
782	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	<i>International Conference on Learning Representa-</i>	838
783	Tran-Johnson, et al. 2022. Language models (mostly)	<i>tions</i> .	839
784	know what they know. <i>CoRR</i> .		
785	Amir Hosein Khasahmadi, Kaveh Hassani, Parsa	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023.	840
786	Moradi, Leo Lee, and Quaid Morris. 2020. Memory-	SelfCheckGPT: Zero-resource black-box hallucina-	841
787	based graph networks. In <i>International Conference</i>	tion detection for generative large language models .	842
788	<i>on Learning Representations</i> .	In <i>Proceedings of the 2023 Conference on Empiri-</i>	843
789	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	<i>cal Methods in Natural Language Processing</i> , pages	844
790	Semantic uncertainty: Linguistic invariances for un-	9004–9017, Singapore. Association for Computa-	845
791	certainty estimation in natural language generation.	tional Linguistics.	846
792	In <i>The Eleventh International Conference on Learn-</i>	Diego Marcheggiani and Ivan Titov. 2017. Encoding	847
793	<i>ing Representations</i> .	sentences with graph convolutional networks for se-	848
794	Chen Li, Yang Liu, Fei Liu, Lin Zhao, and Fuliang	mantic role labeling . In <i>Proceedings of the 2017</i>	849
795	Weng. 2014. Improving multi-documents summa-	<i>Conference on Empirical Methods in Natural Lan-</i>	850
796	rization by sentence compression based on expanded	<i>guage Processing</i> .	851
797	constituent parse trees . In <i>Proceedings of the 2014</i>	Mahdi Pakdaman Naeini, Gregory Cooper, and Milos	852
798	<i>Conference on Empirical Methods in Natural Lan-</i>	Hauskrecht. 2015. Obtaining well calibrated proba-	853
799	<i>guage Processing (EMNLP)</i> , pages 691–701, Doha,	bilities using bayesian binning. In <i>Proceedings of the</i>	854
800	Qatar. Association for Computational Linguistics.	<i>AAAI conference on artificial intelligence</i> , volume 29.	855

856	Dimitrios P Panagoulas, Maria Virvou, and George A	Francesco Tonolini, Nikolaos Aletras, Jordan Massiah,	911
857	Tsihrintzis. 2024. Evaluating llm-generated multi-	and Gabriella Kazai. 2024. Bayesian prompt ensem-	912
858	modal diagnosis from medical images and symptom	bles: Model uncertainty estimation for black-box	913
859	analysis. <i>arXiv preprint arXiv:2402.01730</i> .	large language models . In <i>Findings of the Association</i>	914
		<i>for Computational Linguistics: ACL 2024</i> , pages	915
860	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	12229–12272, Bangkok, Thailand. Association for	916
861	Jing Zhu. 2002. Bleu: a method for automatic evalu-	Computational Linguistics.	917
862	ation of machine translation. In <i>Proceedings of the</i>		
863	<i>40th annual meeting of the Association for Computa-</i>	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	918
864	<i>tional Linguistics</i> , pages 311–318.	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	919
		Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	920
865	Victor Quach, Adam Fisch, Tal Schuster, Adam Yala,	Bhosale, et al. 2023. Llama 2: Open founda-	921
866	Jae Ho Sohn, Tommi S Jaakkola, and Regina Barzi-	tion and fine-tuned chat models. <i>arXiv preprint</i>	922
867	lay. 2024. Conformal language modeling. In <i>The</i>	<i>arXiv:2307.09288</i> .	923
868	<i>Twelfth International Conference on Learning Repre-</i>		
869	<i>sentations</i> .	Boshi Wang, Xiang Yue, and Huan Sun. 2023. Can	924
		chatgpt defend its belief in truth? evaluating llm	925
870	Siva Reddy, Danqi Chen, and Christopher D. Manning.	reasoning via debate. In <i>The 2023 Conference on</i>	926
871	2019. CoQA: A conversational question answering	<i>Empirical Methods in Natural Language Processing</i> .	927
872	challenge . <i>Transactions of the Association for Com-</i>		
873	<i>putational Linguistics</i> , 7:249–266.	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	928
		Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,	929
874	Parshin Shojaee, Kazem Meidani, Shashank Gupta,	Shaokun Zhang, Jiale Liu, et al. 2024. Autogen: En-	930
875	Amir Barati Farimani, and Chandan K Reddy. 2024.	abling next-gen llm applications via multi-agent con-	931
876	Llm-sr: Scientific equation discovery via program-	versations. In <i>First Conference on Language Model-</i>	932
877	ming with large language models. <i>CoRR</i> .	<i>ing</i> .	933
		Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie	934
878	Zixing Song and Irwin King. 2022. Hierarchical het-	Neiswanger, Ruslan Salakhutdinov, and Louis-	935
879	erogeneous graph attention network for syntax-aware	Philippe Morency. 2022. Uncertainty quantification	936
880	summarization. In <i>Proceedings of the AAAI Con-</i>	with pre-trained language models: A large-scale em-	937
881	<i>ference on Artificial Intelligence</i> , volume 36, pages	pirical analysis . In <i>Findings of the Association for</i>	938
882	11340–11348.	<i>Computational Linguistics: EMNLP 2022</i> , pages	939
		7273–7284, Abu Dhabi, United Arab Emirates. As-	940
883	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	sociation for Computational Linguistics.	941
884	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,		
885	and Tatsunori B. Hashimoto. 2023. Stanford alpaca:	Jiacheng Xu and Greg Durrett. 2019. Neural extractive	942
886	An instruction-following llama model. https://	text summarization with syntactic compression . In	943
887	github.com/tatsu-lab/stanford_alpaca .	<i>Proceedings of the 2019 Conference on Empirical</i>	944
		<i>Methods in Natural Language Processing and the</i>	945
888	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	<i>9th International Joint Conference on Natural Lan-</i>	946
889	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 3292–	947
890	Schalkwyk, Andrew M Dai, Anja Hauth, Katie	3303, Hong Kong, China. Association for Computa-	948
891	Millican, et al. 2023. Gemini: a family of	tional Linguistics.	949
892	highly capable multimodal models. <i>arXiv preprint</i>		
893	<i>arXiv:2312.11805</i> .	Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie	950
		Jegelka. 2019. How powerful are graph neural net-	951
894	Katherine Tian, Eric Mitchell, Allan Zhou, Archit	works? In <i>International Conference on Learning</i>	952
895	Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn,	<i>Representations</i> .	953
896	and Christopher Manning. 2023a. Just ask for cali-		
897	bration: Strategies for eliciting calibrated confidence	Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang	954
898	scores from language models fine-tuned with human	Ren, Will Hamilton, and Jure Leskovec. 2018. Hi-	955
899	feedback . In <i>Proceedings of the 2023 Conference</i>	erarchical graph representation learning with differ-	956
900	<i>on Empirical Methods in Natural Language Process-</i>	entiable pooling. <i>Advances in neural information</i>	957
901	<i>ing</i> , pages 5433–5442, Singapore. Association for	<i>processing systems</i> , 31.	958
902	Computational Linguistics.		
		Hao Yuan and Shuiwang Ji. 2020. Structpool: Struc-	959
903	Katherine Tian, Eric Mitchell, Allan Zhou, Archit	tured graph pooling via conditional random fields. In	960
904	Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn,	<i>Proceedings of the 8th international conference on</i>	961
905	and Christopher D Manning. 2023b. Just ask for cali-	<i>learning representations</i> .	962
906	bration: Strategies for eliciting calibrated confidence		
907	scores from language models fine-tuned with human	Xujiang Zhao, Feng Chen, Shu Hu, and Jin-Hee Cho.	963
908	feedback. In <i>Proceedings of the 2023 Conference</i>	2020. Uncertainty aware semi-supervised learning	964
909	<i>on Empirical Methods in Natural Language Processing</i> ,	on graph data. <i>Advances in Neural Information Pro-</i>	965
910	pages 5433–5442.	<i>cessing Systems</i> , 33:12827–12836.	966

Appendix

A Implementation Details

This section provides an overview of the implementation details of GENUINE.

A.1 Details of Datasets

Here in this section, we provide more details on the datasets.

Question Answering. We use the CoQA (Reddy et al., 2019) and TriviaQA (Joshi et al., 2017) datasets to assess LLMs’ ability to generate responses based on contextual understanding and pre-trained knowledge. Additionally, we include the Finance QA dataset (Taori et al., 2023), which evaluates domain-specific knowledge in financial contexts. Rouge-1 (Lin and Och, 2004) is used as the scoring function, labeling a response y_i as correct if $s(y_i, y_{i,true}) \geq 0.3$.

Machine Translation. We evaluate translation quality using the WMT 2014 dataset (Bojar et al., 2014), with BLEU score (Papineni et al., 2002) as the metric. A response y_i is considered correct if $s(y_i, y_{i,true}) \geq 0.3$.

Summarization. The CNN (Hermann et al., 2015) dataset is used for summarization task, where generated outputs are labeled as correct if they achieve a Rouge-L score of at least 0.35, following (Quach et al., 2024).

A.2 Details of Features

This section provides the mathematical definitions of the features used in our uncertainty estimation framework. A detailed breakdown is presented in Table 3.

Table 3: Features used for the supervised task of uncertainty estimation for LLMs.

Name	Definition
Ent	$H(p_\theta(\cdot \mathbf{x}, y_1, \dots, y_{j-1}))$
Max Ent	$\max_{j \in \{1, \dots, n\}} H(p_\theta(\cdot \mathbf{x}, y_1, \dots, y_{j-1}))$
Min Ent	$\min_{j \in \{1, \dots, n\}} H(p_\theta(\cdot \mathbf{x}, y_1, \dots, y_{j-1}))$
Avg Ent	$\frac{1}{n} \sum_{j=1}^n H(p_\theta(\cdot \mathbf{x}, y_1, \dots, y_{j-1}))$
Std Ent	$\sqrt{\frac{\sum_{j=1}^n (H(p_\theta(\cdot \mathbf{x}, y_1, \dots, y_{j-1})) - \text{Avg Ent})^2}{n-1}}$
Prob	$p_\theta(y_j \mathbf{x}, y_1, \dots, y_{j-1})$
Max Prob	$\max_{j \in \{1, \dots, n\}} p_\theta(y_j \mathbf{x}, y_1, \dots, y_{j-1})$
Min Prob	$\min_{j \in \{1, \dots, n\}} p_\theta(y_j \mathbf{x}, y_1, \dots, y_{j-1})$
Avg Prob	$\frac{1}{n} \sum_{j=1}^n p_\theta(y_j \mathbf{x}, y_1, \dots, y_{j-1})$
Std Prob	$\sqrt{\frac{\sum_{j=1}^n (p_\theta(y_j \mathbf{x}, y_1, \dots, y_{j-1}) - \text{Avg Prob})^2}{n-1}}$

A.3 Prompt Template

We adopt a few-shot prompting strategy, following the approach of (Liu et al., 2024b). Each prompt comprises four components: introduction, examples, question, and answer. The examples are user-defined question-answer pairs structured identically to the target task, ensuring consistency in format. The model receives the formatted template along with the reference question and is prompted to generate an appropriate response. This structured approach helps standardize uncertainty estimation across different tasks.

TriviaQA

Answer the question as following examples. Examples: Q: What star sign is Michael Caine? A: Pisces. Q: Which George invented the Kodak roll-film camera? A: Eastman. Q: ... A: ...
Q: In which decade was Arnold Schwarzenegger born? A: 1950s

CoQA

Reading the passage and answer given questions accordingly. Passage: The Vatican Apostolic Library, more commonly called the Vatican Library or simply the Vat, is the library of the Holy See, located in Vatican City. ... Examples: Q: When was the Vat formally opened? A: It was formally established in 1475. Q: ... A: ...
Q: what was started in 2014? A: a project.

WMT

What is the English translation of the following sentence? Q: Spectaculaire saut en wingsuitäu-dessus de Bogota. A: Spectacular Wingsuit Jump Over Bogota. Q: ... A: ...
Q: Une boîte noire dans votre voiture ? A: A black box in your car?

Finance

Answer the question as following examples. Examples: Q: For a car, what scams can be plotted with 0% financing vs rebate? A: he car deal makes money 3 ways. If you pay in one lump payment. ... Q: ... A: ...
Q: Where should I be investing my money? A: Pay off your debt. As you witnessed, no "investment" % is guaranteed. ...

Finance

What are the highlights in this paragraph?
Examples: Q: LONDON, England (Reuters) – Harry Potter star Daniel Radcliffe gains access to a reported £20 million (\$41.1 million) fortune ... A: Harry Potter star Daniel Radcliffe gets £20M fortune as he turns 18 Monday . . . Q: ... A: ...
Q: Editor’s note: In our Behind the Scenes series, CNN correspondents share ... A: Mentally ill inmates in Miami are housed on the “forgotten floor” ...

B Additional Experiments

In this section, we first assess model calibration performance through ECE, NLL, and Brier score metrics, shown in Fig. 6, Fig. 7, and Fig. 8, respectively, comparing GENUINE’s reliability against baselines. Then, we present additional experimental results evaluating GENUINE across four key dimensions: (1) the impact of dependency parse trees on uncertainty estimation (Section B.2), (2) a parameter analysis to determine the sensitivity of GENUINE to hyperparameter tuning (Section B.4), (3) the impact of LLM parameters on GENUINE’s uncertainty estimation performance (Section B.5), and (4) the impact of training dataset size and noisy labels on GENUINE’s uncertainty estimation performance (Section B.6). All experiments are conducted on a Linux server with 64 AMD EPYC 7313 CPUs and an Nvidia Tesla A100 SXM4 GPU with 80 GB of memory.

B.1 Calibration Performance of GENUINE

Calibration ensures that model confidence aligns with actual correctness, making uncertainty estimation more reliable and interpretable. We assess GENUINE and baseline methods using Expected Calibration Error (ECE), Negative Log-Likelihood (NLL), and Brier score, as shown in Fig. 6, Fig. 7, and Fig. 8.

The ECE results (Fig. 6) reveal that while GENUINE outperforms baselines in WMT, Finance, and CNN datasets, it does not consistently achieve the lowest calibration error in TriviaQA and CoQA. This suggests that token-level methods such as SAR and entropy-based approaches remain competitive in capturing uncertainty effectively for shorter responses. A simple guess that GENUINE underperforms in ECE on TriviaQA is due to the data distribution of the TriviaQA dataset, as shown in

Table 6. The rouge score for the TriviaQA dataset is not smooth enough, which can bring bias when using the ECE metric. The ECE metric measures the performance based on each bin(group). The number of samples in each bin can be imbalanced due to the distribution of the Rouge score. Thus, we introduce two other calibration metrics, NLL and Brier score, which focus on measuring the calibration gap at the individual level. However, in longer text generation tasks, where error accumulation can distort confidence estimates, GENUINE demonstrates superior calibration by leveraging dependency structures to refine uncertainty aggregation. The NLL results (Fig. 7) further reinforce these trends. GENUINE consistently achieves lower NLL across all datasets, indicating that it assigns more accurate probability distributions to correct and incorrect responses compared to baselines. The advantage is particularly pronounced in WMT, Finance, and CNN datasets, where long-form responses make token-level uncertainty estimation less effective. Baselines like A4C and SE, which rely on self-evaluation or direct entropy measures, exhibit significantly higher NLL, suggesting that they struggle to generalize confidence estimates across diverse text lengths and response structures.

The Brier score results (Fig. 8) show that GENUINE achieves competitive performance across all datasets, with particularly strong improvements in WMT, Finance, and CNN datasets, aligning with its NLL performance. The gap between GENUINE and its grey-box and white-box variants indicates that hidden layer representations significantly improve calibration, especially for longer outputs. However, the higher ECE in TriviaQA and CoQA suggests that while structural modeling improves overall uncertainty estimation, it may not always provide the best confidence calibration for shorter text generations, where simpler token-wise approaches remain effective.

These results highlight that GENUINE excels in modeling uncertainty for long-form text but is less dominant in short-response tasks, where entropy-based methods can still provide competitive calibration. The findings reinforce the need for task-specific uncertainty estimation strategies, where dependency-aware modeling is particularly beneficial for applications involving complex text structures and extended reasoning.

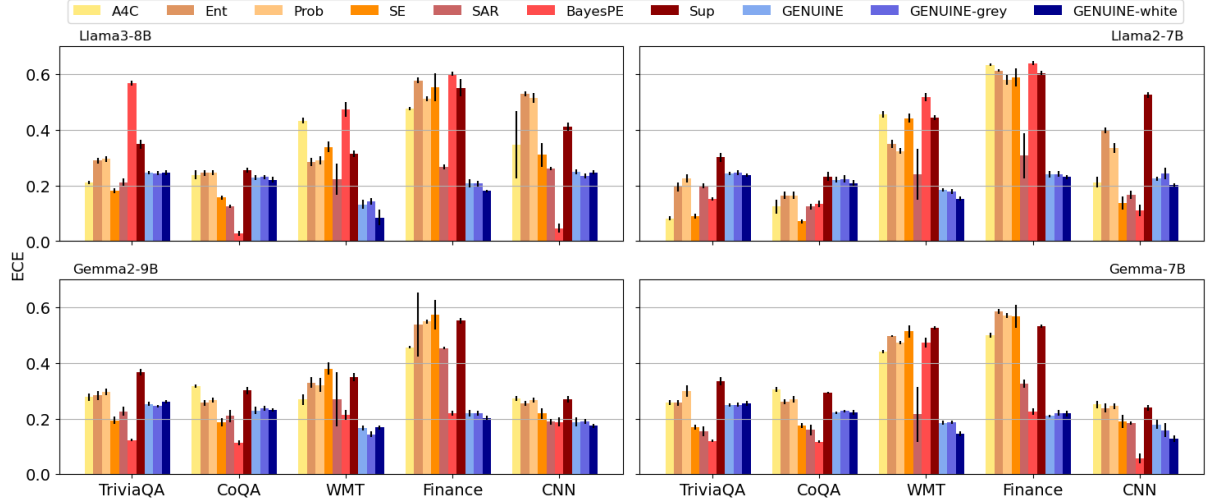


Figure 6: Comparison of ECE on five datasets, four LLMs, and seven baselines. Error bars denote variance over five runs.

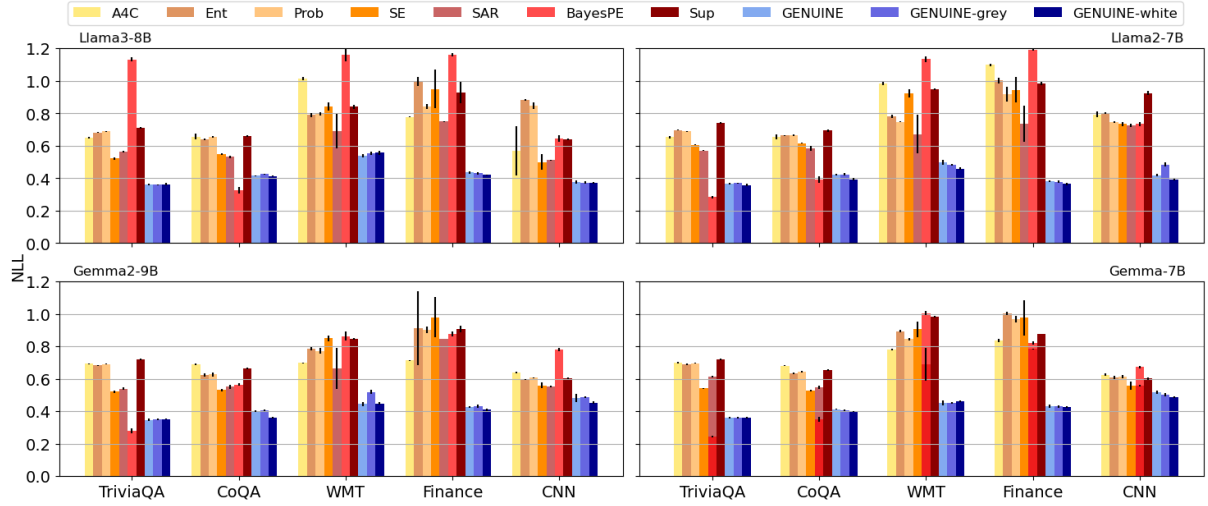


Figure 7: Comparison of NLL on five datasets, four LLMs, and seven baselines. Error bars denote variance over five runs.

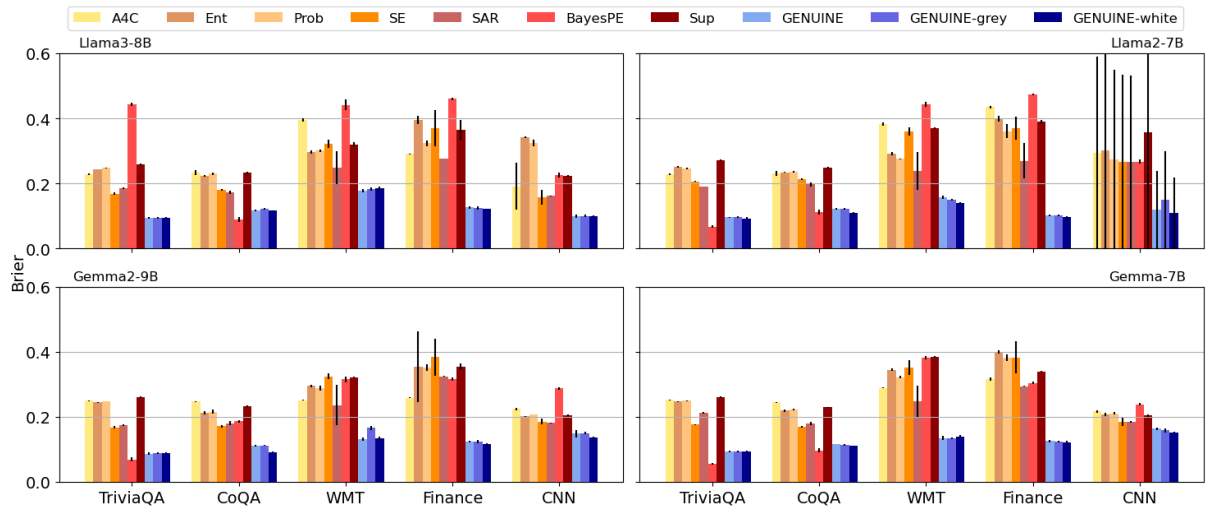


Figure 8: Comparison of Brier scores on five datasets, four LLMs, and seven baselines. Error bars denote variance over five runs.

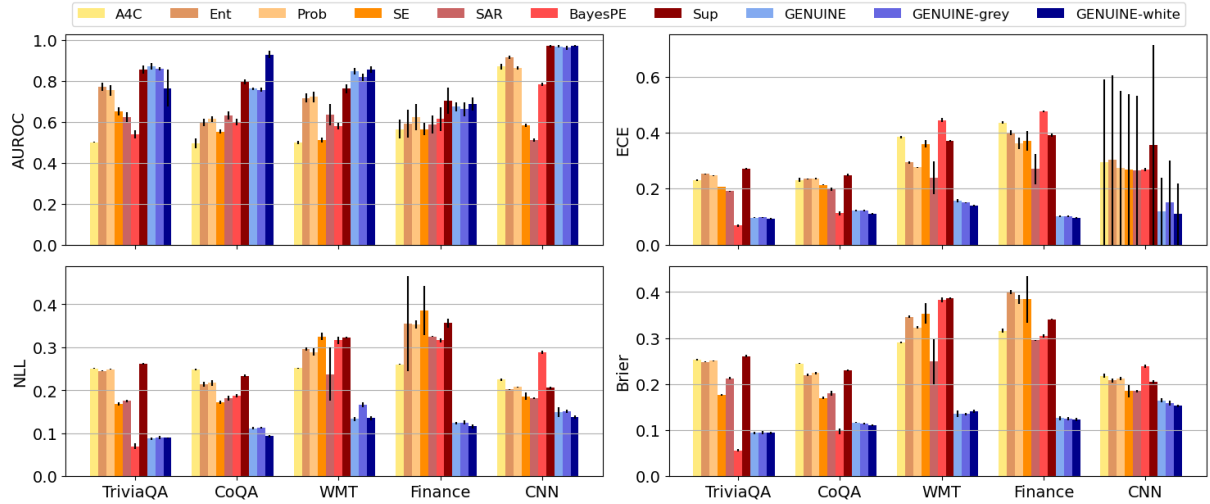


Figure 9: Experimental results on five datasets and seven baseline models on Llama2-13B model. Error bars denote variance over five runs.

Table 4: Comparison of different graph structures for uncertainty estimation on TriviaQA. NTG refers to the next-token graph utilizing both white-box and grey-box features, while DPT represents the dependency parse tree graph with the same feature set. NTG w/ grey and DPT w/ grey denote the respective graphs using only grey-box features, whereas NTG w/ white and DPT w/ white correspond to configurations using only white-box features.(\uparrow means the higher the better, \downarrow means the lower the better)

Graphs	Llama3-8B				Gemma2-9B			
	AUROC \uparrow	ECE \downarrow	NLL \downarrow	Brier \downarrow	AUROC \uparrow	ECE \downarrow	NLL \downarrow	Brier \downarrow
NTG	0.885 \pm 0.048	0.264 \pm 0.040	0.437 \pm 0.133	0.130 \pm 0.062	0.846 \pm 0.088	0.312 \pm 0.082	0.442 \pm 0.122	0.131 \pm 0.056
DPT	0.894 \pm 0.032	0.246 \pm 0.007	0.362 \pm 0.005	0.094 \pm 0.002	0.905 \pm 0.041	0.248 \pm 0.009	0.356 \pm 0.004	0.092 \pm 0.002
NTG w/ grey	0.897 \pm 0.039	0.245 \pm 0.007	0.363 \pm 0.007	0.095 \pm 0.003	0.914 \pm 0.041	0.251 \pm 0.006	0.354 \pm 0.006	0.091 \pm 0.003
DPT w/ grey	0.903 \pm 0.025	0.244 \pm 0.008	0.360 \pm 0.003	0.094 \pm 0.002	0.922 \pm 0.021	0.245 \pm 0.005	0.352 \pm 0.006	0.090 \pm 0.003
NTG w/ white	0.795 \pm 0.049	0.249 \pm 0.010	0.364 \pm 0.007	0.095 \pm 0.003	0.960 \pm 0.019	0.261 \pm 0.009	0.357 \pm 0.006	0.092 \pm 0.003
DPT w/ white	0.809 \pm 0.044	0.246 \pm 0.009	0.362 \pm 0.007	0.094 \pm 0.003	0.970 \pm 0.010	0.261 \pm 0.006	0.353 \pm 0.003	0.090 \pm 0.001

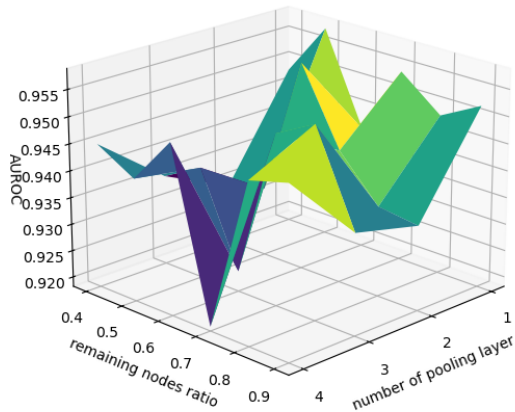


Figure 10: Parameter analysis test on number of pooling layers and remaining nodes ratio for each pooling layer

Table 5: Graph Statistics. Here # Node denotes the average node number and Density denotes the average edge density.

Datasets	Llama3-8B		Llama2-7B	
	# Node	Density	# Node	Density
TriviaQA	3.86	0.56	3.77	0.58
CoQA	5.60	0.47	5.59	0.50
WMT	24.01	0.11	21.75	0.13
Finance	46.46	0.05	21.70	0.15
CNN	61.21	0.04	87.98	0.11

Datasets	Gemma2-9B		Gemma-7B	
	# Node	Density	# Node	Density
TriviaQA	3.83	0.56	3.84	0.56
CoQA	5.28	0.47	5.19	0.48
WMT	23.65	0.12	27.14	0.10
Finance	43.61	0.05	42.92	0.06
CNN	175.33	0.01	162.96	0.01

B.2 Graph Structure and Uncertainty Estimation

Understanding the impact of graph structure on uncertainty estimation is essential for refining confidence assessment in LLM-generated responses. This section evaluates the effectiveness of dependency parse trees and analyzes graph structure variations across datasets and LLMs, using results from

Table 6: Distribution of ROUGE score in TriviaQA. The Bin ID indicates the bin index, the Bin Start indicates the start ROUGE score of the selected bin index, and the Bin End indicates the end ROUGE score of the selected bin index. The Density indicates the ratio of samples in this bin index over the total number of samples.

Bin ID	Bin Start	Bin End	Density
0	0.125	0.213	0.034
1	0.213	0.300	0.383
2	0.300	0.388	0.634
3	0.388	0.475	0.446
4	0.475	0.563	1.131
5	0.563	0.650	0.091
6	0.650	0.738	0.291
7	0.738	0.825	0.520
8	0.825	0.913	0.320
9	0.913	1.000	7.577

Table 4 and Table 5.

Dependency Parse Trees vs. Next-Token Graphs.

To assess the impact of different graph structures, we compare the dependency parse tree (DPT) against the next-token graph (NTG), where edges only connect adjacent words in a sentence. The results in Table 4 clearly demonstrate that DPT-based graphs consistently outperform NTG-based graphs across all evaluation metrics, reinforcing the importance of semantic structure in uncertainty estimation.

For Llama3-8B, DPT achieves an AUROC of 0.894, improving over NTG (0.885), while also achieving lower ECE (0.246 vs. 0.264), NLL (0.362 vs. 0.437), and Brier score (0.094 vs. 0.130). Similar trends hold for Gemma2-9B, where DPT significantly outperforms NTG with an AUROC improvement of nearly 6% (0.905 vs. 0.846) and lower calibration errors. These results confirm that structural relationships encoded in dependency graphs improve uncertainty estimation, providing richer contextual information than simple word adjacency models.

When comparing grey-box vs. white-box features, we observe that DPT consistently performs better than NTG in both settings. For instance, DPT w/ grey achieves an AUROC of 0.903 for Llama3-8B, outperforming NTG w/ grey (0.897) while maintaining better calibration across ECE, NLL, and Brier scores. The trend holds for white-box features, where DPT w/ white achieves 0.809 AUROC vs. 0.795 for NTG w/ white, showing that dependency parsing enhances uncertainty modeling even when using only hidden-layer embeddings.

These findings suggest that semantic-aware uncertainty estimation is essential, especially for

longer text sequences where sequential token dependencies alone fail to capture structural nuances. By modeling hierarchical relations, DPT-based uncertainty estimation improves both reliability and calibration, making it particularly useful for structured prediction tasks.

Graph Variations Across Datasets and LLMs. Beyond structural differences, graph complexity varies significantly across datasets and LLM architectures, as shown in Table 5. We observe several key trends.

First, dataset complexity impacts graph structure. TriviaQA produces the shortest outputs, leading to small graphs with an average of 3.8 nodes, while CNN generates significantly longer responses, resulting in much larger graphs (61.2 nodes for Llama3-8B, 175.3 for Gemma2-9B). This confirms that longer text generations create more intricate dependency structures, further reinforcing why graph-based uncertainty estimation is particularly beneficial for longer responses.

Second, LLM architectures influence graph statistics. While Llama models tend to produce slightly longer responses than Gemma models in shorter datasets like TriviaQA and CoQA, this trend reverses in long-form datasets such as CNN, where Gemma models generate significantly longer outputs than Llama models (e.g., 175.3 nodes vs. 61.2 nodes in CNN for Gemma2-9B and Llama3-8B, respectively). This suggests that some LLM families prioritize brevity while others favor more detailed responses, impacting uncertainty estimation requirements.

Lastly, graph density plays a role in structural complexity. Datasets with shorter outputs (TriviaQA, CoQA) tend to have higher edge density, while longer outputs (CNN, Finance) exhibit lower density, indicating that dependency structures become more sparse as response length increases. This suggests that uncertainty estimation models should be designed to handle both dense, local dependencies and sparse, long-range relationships effectively.

Impact on Uncertainty Estimation Performance: The trends in graph statistics correlate directly with AUROC improvements in Fig. 4, showing that graph-based uncertainty estimation is particularly beneficial for longer text. The WMT dataset, for example, shows substantial AUROC gains when using graph structures, emphasizing that graph-based methods provide the most value in tasks requiring extended reasoning and structured generation.

Table 7: Ablation study of fusion process on TriviaQA(\uparrow means the higher the better)

Methods	Llama3-8B	Llama2-7B
	AUROC \uparrow	AUROC \uparrow
GENUINE w/o fusion & graph	0.789 \pm 0.031	0.835 \pm 0.005
GENUINE w/o fusion	0.809 \pm 0.096	0.843 \pm 0.011
GENUINE	0.894 \pm 0.032	0.860 \pm 0.027
Methods	Gemma2-9B	Gemma-7B
	AUROC \uparrow	AUROC \uparrow
GENUINE w/o fusion & graph	0.956 \pm 0.002	0.853 \pm 0.033
GENUINE w/o fusion	0.963 \pm 0.015	0.900 \pm 0.037
GENUINE	0.969 \pm 0.009	0.917 \pm 0.047

Overall, these findings confirm that dependency parsing enhances uncertainty estimation by providing hierarchical token relationships, making it particularly valuable for long-form generation, structured prediction, and document-level tasks. The graph structure directly influences uncertainty estimation effectiveness, reinforcing the need for adaptive modeling strategies based on dataset and model characteristics.

B.3 Ablation Study

To further prove the effectiveness of the graph structure and the fused assignment matrix, we offer more ablation experiments on the TriviaQA dataset using Llama2-7B and Gemma-7B. As shown in Table 7, the fusion process (Fig. 3) improves AUROC by 2.02% for Llama2-7B and 1.89% for Gemma-7B, the graph structure process improves AUROC by 1.0% for Llama2-7B and 5.5% for Gemma-7B. These results demonstrate that the graph structure and the fusion strategy effectively integrate structural and semantic uncertainty signals, enabling more robust uncertainty propagation across tokens. In contrast, methods w/o a graph structure and fusion strategy fail to capture meaningful relationships between uncertainty features, leading to sub-optimal performance. The consistent improvement across models highlights the importance of structured features and the fusion process in uncertainty estimation. By jointly optimizing structural and semantic representations, GENUINE enhances both robustness and interpretability, making it well-suited for uncertainty-aware applications.

B.4 Parameter Sensitivity

Understanding the impact of hyperparameters on GENUINE’s performance is essential for optimizing uncertainty estimation while ensuring efficiency. We evaluate two key parameters: the number of pooling layers (ranging from 1 to 4) and the remaining node ratio at each pooling step. The

results, shown in Fig. 10, reveal important trends that highlight GENUINE’s robustness and adaptability.

The results indicate that AUROC remains high with fewer pooling layers, suggesting that a deep hierarchy is not necessary for effective uncertainty estimation. As the number of pooling layers increases, performance fluctuates, indicating that excessive pooling may lead to loss of critical structural information, reducing the model’s ability to capture meaningful uncertainty signals. This trend suggests that GENUINE achieves optimal results with a moderate number of pooling layers, avoiding unnecessary complexity while maintaining strong predictive performance.

Additionally, the remaining node ratio plays a crucial role in uncertainty estimation. The model may struggle with redundant information when too many nodes are retained, leading to slightly lower AUROC. However, when the number of retained nodes is optimized, performance improves, reinforcing the idea that removing less informative nodes enhances uncertainty representation. Interestingly, when the remaining ratio is lower, but the number of pooling layers is set appropriately, AUROC reaches peak performance, highlighting the benefits of structured feature reduction in refining uncertainty quantification.

Overall, these findings demonstrate that GENUINE is robust to hyperparameter choices, requiring minimal tuning to achieve strong performance. The ability to maintain high AUROC across a range of configurations suggests that GENUINE can be easily applied to various tasks and LLMs without extensive parameter optimization, making it highly adaptable for real-world deployment.

B.5 Impact of LLM Parameters

Understanding how LLM architecture and scale affect uncertainty estimation is crucial for assessing the generalizability of GENUINE. We compare the performance of Llama2-13B (Fig. 9) against Llama3-8B and Llama2-7B, analyzing its effectiveness across AUROC, calibration metrics (ECE, NLL, and Brier scores), and overall robustness.

Uncertainty Estimation Across LLM Variants.

Llama2-13B achieves strong AUROC performance across all datasets, often matching or surpassing Llama3-8B and Llama2-7B. The improvements are particularly evident in WMT, Finance, and CNN datasets, where Llama2-13B consistently outperforms its smaller counterparts. This suggests that

larger models benefit from enhanced representation learning, leading to more stable and accurate uncertainty estimation in complex, long-form text generation tasks. However, in TriviaQA and CoQA, the AUROC gains are marginal, indicating that the advantages of increased model size are less pronounced for shorter responses.

Calibration Trends: ECE, NLL, and Brier Score Analysis. One notable observation is that GENUINE outperforms baselines in ECE for TriviaQA and CoQA on Llama2-13B, whereas this trend is not observed in Llama3-8B and Llama2-7B. This suggests that larger models may allow GENUINE to better align confidence scores with correctness probabilities in short-response tasks, where previous versions struggled to outperform entropy-based baselines. The ECE results (Fig. 6) further confirm that in WMT, Finance, and CNN, Llama2-13B achieves lower calibration errors, highlighting its ability to generate better-aligned confidence estimates for longer outputs.

The NLL and Brier score results (Fig. 7 and Fig. 8) reinforce these findings. Llama2-13B consistently achieves lower NLL and Brier scores across datasets, particularly in WMT, Finance, and CNN, where uncertainty estimation benefits from structured confidence propagation. This suggests that larger models improve AUROC and provide better-calibrated uncertainty estimates, making them well-suited for tasks requiring complex reasoning and structured text.

The results indicate that larger models significantly enhance both uncertainty estimation and confidence calibration, particularly in short-response tasks like TriviaQA and CoQA, where GENUINE surpasses entropy-based baselines in ECE for the first time. This suggests that model size can influence calibration effectiveness differently across datasets, with larger architectures improving both long-form uncertainty quantification and short-text confidence alignment. Future research should explore adaptive calibration strategies tailored to different response lengths, ensuring that LLMs remain reliable across diverse NLP applications.

Overall, these findings reinforce that GENUINE scales effectively across different LLM architectures, maintaining robust uncertainty estimation and calibration performance while highlighting areas where model size influences uncertainty quantification.

Table 8: The impact of training dataset size on GENUINE performance. More training data results in higher AUROC, but no significant decrease of AUROC when using at least 20% of the training data. Training dataset size does not have much influence on the calibration metrics(\uparrow means the higher the better, \downarrow means the lower the better)

training size	AUROC \uparrow	ECE \downarrow
0.1	0.813 \pm 0.059	0.239 \pm 0.009
0.2	0.873 \pm 0.020	0.244 \pm 0.007
0.3	0.883 \pm 0.012	0.241 \pm 0.011
0.4	0.854 \pm 0.056	0.241 \pm 0.007
0.5	0.874 \pm 0.015	0.243 \pm 0.010
0.6	0.894 \pm 0.032	0.246 \pm 0.007
training size	NLL \downarrow	Brier \downarrow
0.1	0.361 \pm 0.006	0.094 \pm 0.003
0.2	0.363 \pm 0.006	0.095 \pm 0.003
0.3	0.363 \pm 0.001	0.095 \pm 0.000
0.4	0.360 \pm 0.005	0.094 \pm 0.002
0.5	0.362 \pm 0.004	0.094 \pm 0.002
0.6	0.362 \pm 0.005	0.094 \pm 0.002

B.6 Robustness Test on Training Dataset Size and Noisy Labels

In real-world scenarios, uncertainty estimation models often face limited training data and noisy labels, which can affect performance. To evaluate the robustness of GENUINE under such conditions, we conduct experiments using the Llama3-8B model on the TriviaQA dataset. Table 8 shows how varying training set sizes impact performance, while Table 2 examines the effect of label noise. For the latter, we randomly corrupt a portion of training labels (as specified by the noise ratio) and assess performance on the clean test set. These experiments demonstrate GENUINE’s resilience to data scarcity and label noise, highlighting its applicability in real-world settings.

The results shown in Table 8 indicate that the number of training samples does influence GENUINE’s performance, especially when using only 10% of the training data, the AUROC drops 9.1% compared to the model using 60% training data. However, when the training samples take between 20% and 50% of the whole samples, the performances remain relatively stable. Another observation is that the training dataset size does not have much influence on the ECE, NLL, and Brier score.

From the results in Table 2, we find that the noisy labels have a negative influence on the models’ performance in general. However, GENUINE remains robust when 0.1% of the training samples are polluted. As the noise ratio increases, the AUROC drops significantly, as well as the ECE, NLL, and Brier score. We can conclude that, unlike training dataset size, which has little impact on the cali-

bration metrics, the noise ratio influences not only
the AUROC but also the calibration results.