Unconditional Truthfulness: Learning Unconditional Uncertainty of Large Language Models

Artem Vazhentsev^{2,3} Ekaterina Fadeeva⁴ Rui Xing^{1,5}
Gleb Kuzmin^{6,7} Ivan Lazichny³ Alexander Panchenko^{2,3} Preslav Nakov¹
Timothy Baldwin^{1,5} Maxim Panov¹ Artem Shelmanov¹

¹MBZUAI ²Center for Artificial Intelligence ³Computational Semantics Group ⁴ETH Zürich ⁵The University of Melbourne ⁶Weakly-Supervised NLP Group ⁷Laboratory for Analysis and Controllable Text Generation Technologies RAS {vazhentsev, kuzmin, panchenko}@airi.net ekaterina.fadeeva@inf.ethz.ch {rui.xing, preslav.nakov, timothy.baldwin, maxim.panov, artem.shelmanov}@mbzuai.ac.ae

Abstract

Uncertainty quantification (UQ) has emerged as a promising approach for detecting hallucinations and low-quality output of Large Language Models (LLMs). However, obtaining proper uncertainty scores is complicated by the conditional dependency between the generation steps of an autoregressive LLM because it is hard to model it explicitly. Here, we propose to learn this dependency from attention-based features. In particular, we train a regression model that leverages LLM attention maps, probabilities on the current generation step, and recurrently computed uncertainty scores from previously generated tokens. To incorporate the recurrent features, we also suggest a two-staged training procedure. Our experimental evaluation on ten datasets and three LLMs shows that the proposed method is highly effective for selective generation, achieving substantial improvements over rivaling unsupervised and supervised approaches.¹

1 Introduction

Uncertainty quantification (UQ: Gal and Ghahramani (2016); Fadeeva et al. (2023); Baan et al. (2023); Geng et al. (2024)) is of growing interest in the Natural Language Processing (NLP) community for dealing with hallucinations (Fadeeva et al., 2024) and low-quality generations (Malinin and Gales, 2021) in Large Language Models (LLMs) in an efficient manner. For example, high uncertainty could serve as an indicator that the LLM generation should be discarded as potentially harmful or misleading. This approach is known in the literature as selective generation (Baan et al., 2023).

There are many approaches for detecting hallucinations and low-quality outputs of LLMs (Ji et al., 2023; Min et al., 2023; Chen et al., 2023). However, many of them leverage external knowledge

¹https://github.com/mbzuai-nlp/
llm-tad-uncertainty

sources or a second LLM. Knowledge sources are generally patchy in coverage, while censoring the outputs of a small LLM using a bigger one has a high computational cost and is impractical. We argue that LLMs inherently contain information about the limitations of their own knowledge, and that there should be an efficient way to access this information, which can enable LLM-based applications that are both safe and practical.

While a rich body of UQ techniques has been developed for general text classification and regression tasks (Zhang et al., 2019; He et al., 2020; Xin et al., 2021; Wang et al., 2022; Vazhentsev et al., 2022, 2023; He et al., 2024a), applying UQ to text generation is considerably more challenging. A key difficulty arises from the fact that LLMs generate sequences token by token, making multiple conditionally dependent predictions (Zhang et al., 2023). Since LLMs generate text by conditioning on previously produced tokens, an early hallucination, whether at the beginning or in the middle of a sequence, can propagate, causing subsequent claims to also be incorrect. Crucially, even if the generation of the first claim was highly uncertain, this uncertainty is not taken into account during the subsequent generation process. This means that although the first error may be recognized due to its high uncertainty, all subsequent errors are overlooked because the generation process conditioned on it proceeds with high confidence. Therefore, effective hallucination detection requires accounting for this dependency and propagating uncertainty across generation steps.

In this work, we note that the attention between the generated tokens provides information about the conditional dependency between the generation steps. Previously, there have been several attempts to suggest heuristic approaches to model this dependency (Zhang et al., 2023). We argue that the

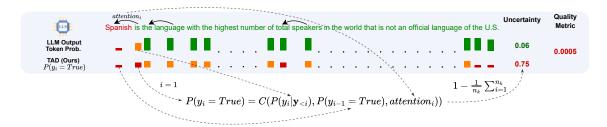


Figure 1: An illustration of the proposed method TAD. The figure shows the generated tokens, the uncertainty scores for the generated sequence, and the probabilities assigned by an LLM and by TAD (represented with bars). The output is generated by LLaMa-3.1 8B for the question What is the language with the highest number of total speakers in the world that is not an official language of the U.S.? The LLM starts by generating the token Spanish that leads to the erroneous answer. The probabilities estimated by the LLM are high for all tokens except for the first one, which makes the uncertainty scores based on raw probabilities misleadingly low. On the contrary, TAD takes into account uncertainty from the previous step using a trainable model $C(\cdot)$ based on attention, resulting in a high overall uncertainty for the generated answer.

particular algorithmic function would be too difficult to engineer, and thus we propose to learn this dependency from data instead.

For this purpose, we generate a training dataset with a target variable, representing the quality score of the generated text according to some ground truth annotation, and train a regression model that leverages LLM attention maps, probabilities on the current generation step, and recurrently computed uncertainty scores from previously generated tokens. To incorporate recurrent features, we suggest a two-staged training procedure where in the second stage, we use scores from the intermediate model obtained in the first training stage. We call the proposed approach *Trainable Attention-based Dependency (TAD)*. Figure 1 illustrates the idea of the method on the real output of an LLM.

The **contributions** of this work are as follows.

- We develop a new data-driven supervised approach to uncertainty quantification that leverages features based on attention maps, probabilities on the current generation step, and recurrently computed uncertainty scores from previously generated tokens.
- We show that both attention and recurrent features are essential for achieving high performance in UQ, and a two-step training procedure is necessary to avoid overfitting.
- We conduct a comprehensive empirical investigation of selective generation, and show
 that the proposed approach outperforms existing unsupervised and supervised UQ methods
 across nine datasets and three LLMs.

2 Related Work

The majority of the methods for UQ for LLMs have been unsupervised, with only a few supervised approaches proposed more recently.

Unsupervised UQ methods. The problem of multiple correct generations was explicitly addressed by Kuhn et al. (2023); Nikitin et al. (2024); Cheng and Vlachos (2024); Zhang et al. (2024) and in a series of black-box generation methods (Lin et al., 2024). The main idea is to sample multiple generations from an LLM, extract semantically equivalent clusters, and analyze the diversity of the generated meanings instead of the surface forms. Chen et al. (2024) proposed evaluating the consistency of the multiple generations in the embedding space using their hidden states. In this category, lexical similarity (Fomicheva et al., 2020) is a very competitive baseline that can be applied to black-box models (without any access to logits or internal model representations). Fadeeva et al. (2024) identified that multiple sources of uncertainty present in the LLM's probability distribution are irrelevant for hallucination detection and proposed a method to mitigate them. Moskvoretskii et al. (2025) addressed the sources of uncertainty in LLMs arising from retrieved context, focusing on determining whether a RAG pipeline should be used, while Belikova et al. (2024) examined how to identify the most suitable context for a given query.

Zhang et al. (2023) and Duan et al. (2024) highlighted that not all tokens should contribute to the uncertainty score, proposing heuristics to select the relevant tokens. Zhang et al. (2023) also modeled the conditional dependencies between the generation steps by penalizing the uncertainty scores based on the uncertainties of the previously generated tokens and the max-pooled attention to the previous tokens.

Overall, most previous work on UQ has not addressed the conditional dependency between the predictions, or has addressed it using heuristics. We argue that the conditional dependency is an important aspect of UQ for text generation tasks, and propose a data-driven approach for dealing with it. We also note that techniques based on sampling multiple answers from LLMs usually introduce prohibitive computational overhead. We argue that for UQ methods to be practical, they should also be computationally efficient.

Supervised UQ methods. Supervised regression-based confidence estimators are well-known for classification problems, primarily from computer vision (Lahlou et al., 2023; Park and Blei, 2024). Their key benefit is computational efficiency.

A handful of papers has applied them to text generation tasks. Lu et al. (2022) proposed training a regression head to predict confidence. They noted that the probability distribution of a language model is poorly calibrated and cannot be used directly to spot low-quality translations. They modified the model architecture and the loss function, restricting this approach to fine-tuning language models only for machine translation and making it unsuitable for general-purpose LLMs. In a similar vein, Azaria and Mitchell (2023) approached the task of UQ by training a multi-layer perceptron (MLP) on the activations of the internal layers of LLMs to classify true vs. false statements. They demonstrated that it outperformed other supervised baselines and few-shot prompting of the LLM itself. However, the reliance on forced decoding limits the real-world applicability to unrestricted generation cases.

Several studies enhanced this method by refining the model architecture and the training procedure. Su et al. (2024) combined the hidden state of the last token with the average hidden state of the sequence, while CH-Wang et al. (2024) introduced a trainable attention layer over token embeddings and used linear regression on top of the MLP's predictions based on embeddings from various layers. He et al. (2024b) proposed to combine multiple deep learning models trained on diverse features extracted from hidden states. Chuang et al. (2024) suggested training the linear classifier using features derived from attention matrices. Vazhentsev

et al. (2025) proposed extracting token embeddings from multiple layers of LLMs, computing density-based scores for each token, and training the linear regression on these features.

Unlike previous methods, we focus on modeling the conditional dependencies between generation steps using attention in a supervised way. The method we propose incorporates recurrently computed uncertainty scores for tokens from previous generation steps, capturing the relationship between the uncertainty of generated tokens. Additionally, our method is flexible as it can be applied at different levels: to the entire text, to a sub-sequence, or to individual tokens. Finally, unlike the method proposed by Chuang et al. (2024), which relies on feature engineering, our method directly utilizes raw attention weights that give access to more information.

3 Problem Background and Key Idea

When an LLM generates a sequence of tokens y_i , it provides us a conditional probability distribution $p(y_i \mid \mathbf{y}_{< i}) = p(y_i \mid \mathbf{x}, \mathbf{y}_{< i})$, where \mathbf{x} is an input prompt and $\mathbf{y}_{< i}$ is a sequence of tokens generated before token y_i . This essentially means that the LLM considers that everything generated so far is correct, which might not be the case. In practice, we would like to somehow propagate the uncertainty from the previous generation steps.

To illustrate the problem, for the sake of simplicity, let us assume that only the uncertainty from the previous tokens is propagated to the current generation step. This assumption can be expressed as follows: $p(y_i \mid \mathbf{y}_{< i}) \simeq p(y_i \mid y_{i-1})$. Let us further consider that we have trained an LLM that generates only tokens that are true ("T") or false ("F"). The probability of the token y_i being T is given by the conditional probability $p(y_i \mid y_{i-1}) = p(y_i = T \mid y_{i-1} = T)$. Assume we already have some tokens y_1, y_2, \ldots, y_n and a prompt \mathbf{x} . At each step, the LLM provides us $p(y_1 = T \mid \mathbf{x}), p(y_2 = T \mid y_1 = T), \ldots, p(y_n = T \mid y_{n-1} = T)$.

These probability distributions are conditionally dependent on the previously generated tokens. However, to estimate the correctness of some token y_i , we need to obtain an *unconditional probability* $p(y_i) = p(y_i = T)$. Let us expand $p(y_i = T)$ according to the law of total probability and express

it using conditional probability:

$$p(y_i = T) = p(y_i = T \mid y_{i-1} = T) \cdot p(y_{i-1} = T) + p(y_i = T \mid y_{i-1} = F) \cdot (1 - p(y_{i-1} = T)).$$

In this formula, $p(y_i = T \mid y_{i-1} = T)$ is what the LLM provides during the current generation step in accordance with the specified assumptions, and $p(y_{i-1} = T)$ is recurrently calculated based on the previous generation step. We still do not know the remaining term: $p(y_i = T \mid y_{i-1} = F)$. This simplistic example shows that in order to obtain an uncertainty estimate suitable for hallucination detection, we cannot rely solely on the probability distribution provided by the LLM, and we also need to model the conditional dependency of the generation steps. It also makes explicit the need for recurrence in token-level uncertainty computation.

Attention weights are commonly used in interpretability methods to illustrate which tokens influenced the model's decision at the current generation step (Zhao et al., 2024; Tufanov et al., 2024; Ferrando and Voita, 2024). However, obtaining a direct expression that would accurately approximate the conditional dependency between the generation steps is challenging. The assumptions in our simplistic example do not hold in real LLMs, and thus the predictions on each step depend on multiple previous tokens in a complicated fashion. We suggest learning this dependency in a supervised way from attention. In particular, we propose a feature set for training token-level unconditional confidence scores C, consisting of the attention weights Att_i , the token probabilities from the LLM on the current step $p(y_i \mid \mathbf{y}_{\leq i})$, and the recurrently calculated confidence scores on the previous steps $C_{< i}$:

$$C(y_i) = C(Att_i, p(y_i \mid \mathbf{y}_{\le i}), \mathbf{C}_{\le i}). \tag{1}$$

4 Trainable Attention-Based Conditional Dependency

We suggest learning unconditional token-level probability estimates from annotated data.

Obtaining targets for learning unconditional probability. In order to obtain the targets $\hat{p}(y_i)$ for the unconditional probability $C(y_i)$ for a generated token $y_i \in \mathbf{y}$ during the training phase, we compute the semantic similarity between the generated answer \mathbf{y} and the ground truth \mathbf{y}^* :

$$\hat{p}(y_i) = \sin(\mathbf{y}, \mathbf{y}^*). \tag{2}$$

For generating the targets, we use task-specific similarity measures, such as Accuracy, COMET (Rei et al., 2020), and AlignScore (Zha et al., 2023).

Generating training data for TAD. We generate the training data for TAD using the original textual training dataset in the following way:

- 1. For the input prompt \mathbf{x}_k , we use an LLM to generate a text $\mathbf{y}_k = y_1 y_2 \dots y_{n_k}$ of some length n_k and token probabilities $p(y_i \mid \mathbf{x}_k, \mathbf{y}_{\leq i})$.
- 2. For the first generated token y_1 in each text, we introduce its unconditional confidence estimate $\hat{p}_k(y_1) = \sin(\mathbf{y}_k, \mathbf{y}_k^*)$ according to Equation (2).
- 3. For each generated token y_i , $i = 2, ..., n_k$ we construct a feature vector z_i^k that depends on N preceding tokens. The feature vector z_i^k includes: the conditional probabilities $p(y_i \mid \mathbf{x}_k, \mathbf{y}_{< i})$ and $p(y_{i-l} \mid \mathbf{x}_k, \mathbf{y}_{< i-l})$, for $l = 1, ..., \min\{N, i - 1\}$; the unconditional probabilities' estimates from the previous steps $\hat{p}_k(y_{i-l})$, and the attention weights $a_{i,i-l}$ from the (i-l)-th token to the i-th token from all layers and heads. If N > i - 1, we pad the feature vector with zeros to ensure they have the same length. During the first training stage, z_i^k includes only the conditional probabilities without other features. Consequently, on the first stage of learning, the unconditional probabilities $\hat{p}_k(y_{i-1})$ are not required. On the subsequent learning stages it is estimated via the function learned on the previous learning stage.

As a result, for each instance in the training dataset and for each iteration of learning, we generate a sequence of target variables $\hat{p}_k(y_i) = \sin(\mathbf{y}_k, \mathbf{y}_k^*)$ and corresponding feature vectors $z_i^k, k = 1, \dots, K, i = 2, \dots, n_k$. We use this dataset to train the model C. The step-by-step procedure for generating training data is presented in Algorithm 1 in Appendix E.

Model for C and its training procedure. The training procedure involves using the estimates of the unconditional probabilities from the previous steps as features. To address this problem, we perform the training procedure twice. In the second stage, we leverage the predictions of the function C trained on the first stage as features. This two-step training approach enables us to leverage the conditional dependency of the current step on the previous ones when computing the uncertainty score.

Our experiments show that it is essential for achieving good performance.

We experiment with two regression models for TAD: linear regression (LinReg) and a multi-layer perceptron (MLP). The hyper-parameters of the regressors are obtained using cross-validation with five folds on the training dataset. We select the optimal values of the hyperparameters based on the best average PRR. The optimal values are used to train the regression model on the full training set. The selected hyper-parameters values for the TAD modules are presented in Appendix C.1.

Inference procedure. During inference, we obtain predictions from the LLM as always, but we also extract features from the attention outputs. For the first generated token y_1 , its unconditional probability is defined as $p(y_1) = p(y_1 \mid \mathbf{x}_k)$. For each subsequent token, the function C computes the predictions recursively, leveraging the attentions, the conditional probabilities, and the unconditional probabilities predicted for the preceding tokens. Finally, to compute the uncertainty of the LLM answer, the token-level scores are aggregated into a sequence-level score:

$$U(\mathbf{y}) = 1 - \frac{1}{n_k} \sum_{i=1}^{n_k} C^k(y_i).$$
 (3)

We experiment with various aggregation approaches in the ablation study.

5 Experiments and Evaluation

5.1 Experimental Setup

For the experimental evaluation, we use the LM-Polygraph framework (Fadeeva et al., 2023). We focus on the task of selective generation (Ren et al., 2023) where we "reject" generated sequences due to low quality based on uncertainty scores. Rejecting means that we do not use the model output, and the corresponding queries are processed differently, e.g., they could be further reprocessed manually.

Evaluation measures. Following previous work on UQ in text generation (Malinin and Gales, 2021; Vashurin et al., 2025; Ielanskyi et al., 2025), we compare UQ methods using the Prediction Rejection Ratio (PRR) metric. PRR quantifies how well an uncertainty score can identify and reject low-quality predictions according to some quality measure. The PRR scores are normalized to the range [0, 1] by linearly scaling the area under the PR curve between the values obtained with random

selection (corresponding to 0) and oracle selection (corresponding to 1). Higher PRR values indicate better quality of the selective generation. Following previous work (Vashurin et al., 2025), we compute PRR only up to a rejection threshold of 50% to ensure its practical applicability. We use Accuracy, COMET (Rei et al., 2020), and AlignScore (Zha et al., 2023; Santilli et al., 2025) as generation quality measures. We also use ROC-AUC of detecting incorrect answers as a supplementary metric, as it is widely adopted in the UQ literature.

Datasets. We consider ten datasets from five text generation tasks: text summarization (TS), machine translation (MT), Question Answering (QA) with long free-form answers, QA with free-form short answers, and multiple-choice QA. A detailed description of all datasets is provided in Appendix D, and the dataset statistics are presented in Table 21.

LLMs. We experiment with three LLMs: LLaMA-3.1 8b (Grattafiori et al., 2024), Gemma-2 9b (Rivière et al., 2024), and Qwen-2.5 7b (Yang et al., 2024). The values of the inference hyperparameters are given in Table 20 in Appendix C.2.

UQ baselines. The set of unsupervised baselines includes Maximum Sequence Probability (MSP), Mean Token Entropy, and Perplexity (Fomicheva et al., 2020), which are considered simple yet strong and robust baselines for selective generation across various tasks (Fadeeva et al., 2023). We also compare our method to unsupervised techniques considered to be state-of-the-art: Lexical Similarity based on ROUGE-L (Fomicheva et al., 2020), black-box methods (DegMat, Eccentricity, EigValLaplacian: Lin et al. (2024)), Semantic Entropy (Kuhn et al., 2023), hallucination detection with a stronger focus (Focus: Zhang et al. (2023)), claim-conditioned probability (CCP: Fadeeva et al. (2024)), Shifting Attention to Relevance (SAR: Duan et al. (2024)), EigenScore (Chen et al., 2024), Semantic Density (Qiu and Miikkulainen, 2024), and long-text uncertainty quantification (LUQ: Zhang et al. (2024)). For samplingbased methods, we generate five samples.

The suite of baselines also includes state-of-the-art supervised methods that use hidden states or attention weights: Factoscope (He et al., 2024b), SAPLMA (Azaria and Mitchell, 2023), MIND (Su et al., 2024), Sheeps (CH-Wang et al., 2024), LookBackLens (Chuang et al., 2024), and SATRMD (Vazhentsev et al., 2025).

UQ Method	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	MedQUAD AlignScore	TruthfulQA AlignScore	CoQA AlignScore	SciQ AlignScore	TriviaQA AlignScore	MMLU Acc.	GSM8k Acc.	Mean PRR	Mean Rank
MSP	.303	.107	.329	.459	.091	.314	.262	.459	.527	.535	.310	.336	10.55
Perplexity	.353	.408	.076	.416	.249	.363	.259	.244	.506	.492	.303	.334	11.36
Mean Token Entropy	.342	.424	.056	.425	.238	.384	.251	.218	.528	.247	.333	.313	11.55
CCP	.339	.065	.338	.363	.038	.142	.210	.351	.562	.446	.306	.287	13.64
Simple Focus	.254	.290	.196	.472	.074	.271	.281	.486	.545	.516	.302	.335	10.45
Focus	.310	.324	.055	.416	.137	.386	.211	.422	.507	.305	.278	.305	12.55
Lexical Similarity Rouge-L	.076	.132	.061	.403	017	.063	.277	.378	.491	.242	.273	.216	16.18
EigenScore	.033	.070	.055	.318	010	.075	.263	.355	.462	.192	.283	.191	18.18
EVL NLI Score entail.	.033	.086	.113	.252	.137	.253	.314	.371	.577	.230	.188	.232	14.55
Ecc. NLI Score entail.	.012	.004	000	.340	.102	.126	.293	.380	.530	.231	.235	.205	16.82
DegMat NLI Score entail.	.031	.089	.113	.285	.146	.253	.316	.429	.583	.239	.203	.244	13.00
Semantic Entropy	.016	.083	.085	.379	.093	.092	.232	.347	.479	.157	.366	.212	16.91
SAR	.052	.166	.049	.435	.107	.145	.297	.439	.552	.275	.320	.258	12.45
LUQ	.137	.210	.147	.224	.101	.212	.303	.394	.570	.249	.158	.246	13.55
Semantic Density	.163	.122	.100	.295	.175	.320	.380	.448	.571	.237	.197	.273	11.82
Factoscope	.292	.064	020	.120	.511	.065	.033	.313	.363	.585	.121	.222	17.45
SAPLMA	.288	.382	.056	.548	.228	.277	002	.399	.399	.456	.358	.308	11.82
MIND	.437	.361	.178	.451	.531	.411	.263	.499	.517	.727	.570	.450	6.36
Sheeps	.510	.466	.380	.509	.501	.349	.423	.552	.594	.723	.604	.510	3.09
LookBackLens	.528	.441	.279	.613	.547	.462	.341	.542	.497	.718	.525	.499	4.45
SATRMD	.494	<u>.495</u>	.248	.475	.424	.448	.333	.581	.561	.704	.528	.481	4.45
TAD	.550	.535	.444	.592	.624	.463	.392	.488	.632	<u>.724</u>	.557	.545	1.82

Table 1: PRR↑ of UQ methods for the Llama-3.1 8b model. Warmer color indicates better results. The best method is in **bold**, the second best is underlined.

UQ Method	Llama-3.1 8b	Gemma-2 9b	Qwen-2.5 7b	Mean Rank
MSP	10.55	10.27	12.91	9.67
Perplexity	11.36	10.91	10.45	8.33
Mean Token Entropy	11.55	10.45	10.55	8.67
CCP	13.64	12.64	14.36	15.67
Simple Focus	10.45	9.45	11.55	7.00
Focus	12.55	10.18	14.55	13.00
Lexical Similarity Rouge-L	16.18	14.82	14.00	16.67
EigenScore	18.18	19.18	16.18	21.33
EVL NLI Score entail.	14.55	14.91	13.09	16.33
Ecc. NLI Score entail.	16.82	17.45	15.45	19.33
DegMat NLI Score entail.	13.00	13.91	12.36	14.00
Semantic Entropy	16.91	16.64	17.18	20.00
SAR	12.45	11.73	12.09	11.33
LUQ	13.55	13.82	12.18	13.33
Semantic Density	11.82	11.64	12.27	11.17
Factoscope	17.45	19.00	17.27	21.33
SAPLMA	11.82	9.36	14.09	10.83
MIND	6.36	7.55	7.36	4.67
Sheeps	3.09	8.09	4.73	3.33
LookBackLens	4.45	4.64	3.55	2.83
SATRMD	4.45	4.00	5.55	3.17
TAD	1.82	2.36	1.27	1.00

Table 2: Mean ranks of UQ methods aggregated over all datasets for each LLM separately (the lower the better). The column *Mean Rank* corresponds to the mean rank of the ranks across all LLMs. The best method is in **bold**, the second best is underlined.

5.2 Main Results

Fine-grained comparison to the baselines. Tables 1, 5 and 6 in Appendix A.1 present the results for LLaMa-3.1 8b, Gemma-2 9b, and Qwen-2.5 7b, respectively.

The results demonstrate that, across all summarization and translation datasets, both LookBack-Lens and TAD outperform state-of-the-art methods by a substantial margin. For Llama, Look-BackLens achieves slightly better results than TAD on the WMT19 dataset, but TAD confidently outperforms LookBackLens on all summarization datasets. For Qwen, TAD consistently achieves the best results on all summarization and translation datasets, while LookBackLens achieves the second-best results.

For QA involving long answers (e.g., MedQUAD, TruthfulQA), TAD demonstrates substantial improvements over the baselines across all considered models. For example, in the experiment with LLaMA-3.1 8b on MedQUAD, TAD outperforms the second-best baseline, LookBack-Lens, by 0.077 of PRR. On the TruthfulQA dataset, TAD achieves an improvement of 0.045 in PRR over the second-best baseline with Gemma. On the GSM8k dataset, TAD consistently demonstrates strong performance and outperforms unsupervised methods, although it performs slightly worse than the Sheeps method.

For QA with short answers (CoQA, SciQ, and TriviaQA), TAD generally exhibits notable improvements over the baseline methods in the majority of cases. The only exception is the case of the SciQ dataset, where LookBackLens is marginally better for LLaMA-3.1 8b and Gemma-2 9b. On TriviaQA, when using the LLaMA-3.1 8b model, TAD outperforms sampling-based methods, while other supervised methods fall behind simple baselines by a margin.

Finally, for MMLU, TAD also notably outperforms state-of-the-art methods for both Gemma-2 9b and Qwen-2.5 7b. However, for LLaMA-3.1 8b, TAD slightly falls behind MIND.

Summarizing, our findings indicate that certain UQ methods, such as LookBackLens, SATRMD, and Sheeps, can achieve top performance in specific experimental settings. However, TAD demonstrates the most consistent and robust performance across all eleven tasks, never ranking below the state-of-the-art unsupervised methods. In contrast,

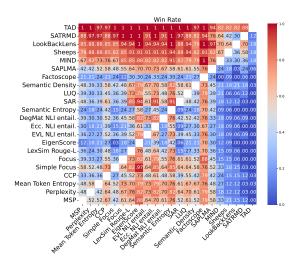


Figure 2: Summary of 33 experimental setups with various models and datasets. Each cell in the diagram presents a fraction of experiments where a method from a row outperforms a method from a column. Warmer colors indicate better results.

other supervised methods occasionally underperform, sometimes even falling below simple baselines such as MSP. Similar patterns are observed in the ROC AUC results reported in Tables 7 to 9 (see Appendix A.2).

Aggregated results. Table 2 presents the mean rank of each method aggregated over all datasets for each model separately. Lower ranks are better. The column *Mean Rank* shows the mean rank of the ranks across all models. Figure 2 additionally summarizes all experimental setups. Each cell presents a win rate for a method from a row compared to a method from a column. The aggregated results emphasize the significance of the performance improvements of the proposed method. Despite some baselines showing good results in particular cases, they are usually quite unstable, resulting in poor overall ranking. TAD demonstrates more robust improvements across multiple tasks and LLMs, making it a better choice overall.

Generalization to out-of-domain datasets. Table 3 compares the results of the supervised methods trained on all QA datasets except for one that represents the out-of-domain dataset for testing. Additionally, Table 10 in Appendix A.3 presents the results when these methods are trained on all QA datasets and tested on the out-of-distribution tasks: summarization and translation. These settings evaluate the out-of-domain generalization capabilities of the supervised techniques for both new domains and new tasks.

UQ Method	CoQA AlignScore	SciQ AlignScore	TriviaQA AlignScore	MMLU Acc.	GSM8k Acc.	Mean PRR
MSP	.262	.459	.527	.535	.310	.419
SAR	.297	.439	.552	.275	.320	.377
Semantic Density	.380	.448	.571	.237	.197	.366
Factoscope	.016	.055	.161	.078	.049	.072
SAPLMA	030	.199	112	089	077	022
MIND	.044	.153	.237	.252	.230	.183
Sheeps	.092	.422	.295	.425	.323	.312
LookBackLens	.079	.365	.304	.422	.166	.267
SATRMD	.247	.349	.469	.205	.311	.316
TAD	.283	.529	<u>.565</u>	<u>.512</u>	.278	.434

Table 3: PRR↑ for Llama-3.1 8b model for various QA tasks for the considered supervised sequence-level methods trained on the general QA dataset. Unsupervised methods are not included as their performance is not dependent of the training data. Warmer colors indicate better results. The best method is in **bold**, and the second best one is underlined.

The results show that all considered supervised methods substantially degrade compared to their indomain performance and, in many cases, underperform the simple MSP baseline. Nevertheless, TAD demonstrates strong out-of-domain performance on the unseen QA datasets, outperforming MSP by 0.015 PRR on average. However, all supervised methods perform significantly worse than the MSP baseline on the OOD tasks, summarization and translation, underscoring their limited adaptability to unseen tasks.

These findings indicate that previous supervised UQ methods are generally effective only for indomain selective generation. However, the TAD method demonstrates the ability to achieve generalization to unseen domains within similar tasks. More details about these experiments are presented in Appendix A.3.

5.3 Ablation Studies

Comparison of features. Table 15 in Appendix A.5 presents the ablation experiment with different features for the TAD regression model. For TAD (probs.), we only use probabilities along with predictions from the preceding tokens $p(y_{i-k} = T)$ for k = 1, ..., N. For TAD (attention), we use attention weights on the N preceding tokens without probabilities. The results show that TAD (probs.) provides meaningful but usually lower performance. TAD (attention) demonstrates substantial improvements, underscoring the importance of using the attentions in the TAD method. Finally, TAD (attention+probs.), which combines all of attention weights, probabilities, and uncertainty scores from previous steps, achieves slight but consistent performance gains. This indicates

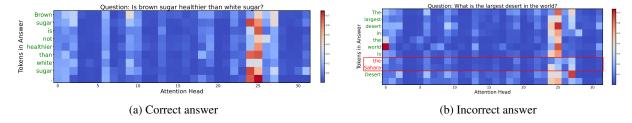


Figure 3: Comparison of the attention weights of Llama-3.1 8B to the last preceding token for each generated token for correct and incorrect answers to input questions from the TruthfulQA dataset. The y-axis shows the generated tokens, and the x-axis represents the attention heads in the 30th layer. Warmer colors indicate higher attention values. In the incorrect answer (Figure b), the model hallucinates the factually incorrect tokens *The Sahara* (the correct answer is Antarctica). Notably, while the 25th attention head consistently assigns high attention to the preceding token in both outputs, this attention noticeably drops for the hallucinated tokens *The Sahara*. This decrease in attention could serve as a valuable signal for a hallucination detector in the TAD method.

the benefit of recurrence during the computation of uncertainty scores.

Impact of the token-level training procedure. Table 14 in Appendix A.5 presents an ablation study comparing different training procedures for the regression model in the TAD method. We compare the original TAD against *TAD (Sequence-level)*, which uses a two-layer MLP with averaging of the hidden features between layers, followed by a linear layer for direct sequence-level uncertainty prediction. The results demonstrate that while *TAD (Sequence-level)* performs competitively: the original TAD method surpasses it by 0.027 PRR on average, with the largest improvement of 0.124 PRR on MedQUAD. These findings highlight the effectiveness of the token-level training procedure with recurrent features in TAD.

Impact of the two-step training procedure. Table 16 in Appendix A.5 presents the ablation experiment comparing one-step vs. two-step training procedures for the TAD method. The results show that the two-step procedure is essential for training a well-performing recurrent model.

Regression models and aggregation approaches. Detailed results with various regression models and aggregation approaches are presented in Table 12. The optimal values of the hyper-parameters of TAD for all experimental setups are presented in Tables 17 to 19 in Appendix C.1 for LLaMA-3.1 8b, Gemma-2 9b, and Qwen-2.5 7b, respectively.

We compared two strategies for aggregating the token-level TAD scores: (*i*) the mean of the scores; and (*ii*) the sum of the log scores, inspired by perplexity. For the majority of the considered settings, the mean of the probabilities yielded the best results. However, for QA with short answers, the sum

of the log probabilities performed slightly better.

We can see that the difference between MLP and LinReg is minimal. On average, TAD with LinReg outperforms TAD with MLP by 0.013 PRR. Therefore, for simplicity, we use LinReg as a regression method for TAD.

Impact of the number of previous tokens. Table 13 in Appendix A.5 presents results with different numbers of preceding tokens used in TAD. It shows that using ten preceding tokens generally yields better performance compared to using only 1–2 tokens across all datasets, except for XSum.

Impact of the attention layers. Figure 4 in Appendix A.5 presents the normalized average weights of linear regression for different attention layers in the TAD method. We can see similar patterns across various tasks, revealing that the most important layers are typically the middle ones, which is consistent with observations in previous work (Azaria and Mitchell, 2023; Chen et al., 2024). Additionally, we note that for the majority of the tasks, the first and last attention layers play a crucial role.

Replacing attention weights with interpretability features. Table 11 in Appendix A.4 shows the results, where we investigate interpretability features from Layer Integrated Gradients (LIG: Sundararajan et al. (2017)) as a measure of conditional dependency between generation steps. We compare the original TAD method with two variants: TAD (LIG), which replaces attention weights with LIG features, and TAD (MIX), which concatenates LIG features with the raw attention weights. LIG features perform comparably to attention, but their inclusion does not enhance TAD performance.

5.4 Analysis of Attention Maps

To better understand the mechanisms behind the state-of-the-art performance of the TAD method, we conducted an analysis of the attention maps used as features in the hallucination detector for Llama-3.1 8B. Figure 3 illustrates the attention weights to the last preceding token of each generated token, for both correct and incorrect answers to input questions from the TruthfulQA dataset. We focus exclusively on attention weights from the current token to its immediate predecessor. We chose this focus because Table 13 in Appendix A.5 shows that relying solely on attention to the previous token still achieves strong performance.

This analysis reveals several key patterns in the attention weights that TAD may leverage when making its predictions. First, we observe that there exists a small number of attention heads that usually pay high attention to the previous token. In our example, the highest average attention across all generated tokens was expressed in the 30th layer by the 25th head. Second, during correct generation, all tokens assign high attention to previous tokens in these heads, whereas in hallucinated outputs, this attention becomes blurred.

For instance, in Figure 3a, all tokens are correct, resulting in consistently high attention from the 25th head. In contrast, in Figure 3b, the model hallucinates the factually incorrect token *The Sahara* (the correct answer is Antarctica), and the attention to this token drops noticeably. This decrease in attention provides a valuable signal for detecting hallucinations in the TAD method.

5.5 Computational Efficiency

To demonstrate the computational efficiency of TAD, we compare its runtime to other UQ methods. We use a single 80GB H100 GPU, as detailed in Table 1. The inference is implemented as a single-batch model call for all tokens in the output.

Table 4 presents the average runtime per text instance for each UQ method, along with the percentage overhead over the standard LLM inference with MSP. As we can see, many state-of-the-art UQ methods (such as DegMat, Lexical Similarity, Semantic Entropy, and SAR) introduce huge computational overhead (400–600%) because they need to perform sampling from the LLM multiple times. In contrast, all supervised methods introduce minimal overhead. In particular, TAD introduces only 5% overhead, which makes it a highly practical and

UQ Method	Runtime per batch	Overhead
MSP	$1.30_{\pm 0.62}$	-
DegMat NLI Score Entail.	$6.86_{\pm 2.28}$	430 %
Lexical Similarity ROUGE-L	$6.72{\scriptstyle\pm2.24}$	420%
Semantic Entropy	$6.86{\scriptstyle\pm2.28}$	430%
SAR	$8.83_{\pm 2.94}$	580%
Factoscope	3.30±2.13	150%
SAPLMA	$1.30_{\pm 0.62}$	0.06%
MIND	$1.30_{\pm 0.62}$	0.10%
Sheeps	$1.50_{\pm 0.97}$	15%
LookBackLens	$1.30_{\pm 0.62}$	0.08%
SATRMD	$1.39_{\pm 0.67}$	8%
TAD	1.37±0.68	5%

Table 4: Evaluation of the inference runtime of UQ methods measured on all test instances from all datasets with predictions from Llama-3.1 8b. The best results are in **bold**, and the second best results are <u>underlined</u>.

efficient choice for uncertainty quantification.

6 Conclusion and Future Work

We have presented a new uncertainty quantification method based on learning conditional dependencies between the predictions made on multiple generation steps. The method relies on attention to construct features for learning this functional dependency and leverages this dependency to alter the uncertainty of the subsequent generation steps. This yields improved results in selective generation tasks, especially when the LLM output is long. Our experimental study shows that TAD usually outperforms other state-of-the-art UQ methods (such as SAR) resulting in the best overall performance across three LLMs and nine datasets. Contrary to other supervised methods, TAD also shows cross-domain generalization. Our method requires only minimal computational overhead due to the simplicity of the underlying linear regression model, making it a practical choice for LLM-based applications.

In future work, we aim to apply the suggested method to UQ of retrieval-augmented LLMs. TAD potentially could be used to assess the credibility of a retrieved piece of textual evidence.

Limitations

The proposed approach is supervised and thus benefits from task-specific training data. We evaluate our method on out-of-domain data to explore its generalization. Despite expected variations in per-

formance, the proposed method achieves promising results on unseen out-of-domain data when trained on the related source domain. Overall, the method can be used in out-of-domain settings, while caution should be exercised when training on significantly different domains.

Our experiments were conducted using 7–9B parameter models, due to limitations in our available computational resources. Nevertheless, given the similar architectures and training procedures across model scales, we believe that the proposed method can be effectively applied to larger-scale LLMs.

Ethical Considerations

In our work, we considered open-weights LLMs and datasets not aimed at harmful content. However, LLMs may generate potentially damaging texts for various groups of people. Uncertainty quantification techniques can help create more reliable use of neural networks. Moreover, they can be applied to detecting harmful generation, but this is not the target of this paper.

Moreover, despite our proposed method demonstrating sizable performance improvements, it can still mistakenly highlight correct and innocuous generated text with high uncertainty in some cases. Thus, as with other uncertainty quantification methods, it has limited applicability.

In the writing of this paper, we used writing assistants to ensure grammatical accuracy.

Acknowledgments

We sincerely thank the anonymous reviewers for their insightful comments and suggestions, which have greatly strengthened this paper

References

- Asma Ben Abacha and Dina Demner-Fushman. 2019. A question-entailment approach to question answering. *BMC Bioinform.*, 20(1):511:1–511:23.
- Amos Azaria and Tom Mitchell. 2023. The internal state of an LLM knows when it's lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, Singapore. Association for Computational Linguistics.
- Joris Baan, Nico Daheim, Evgenia Ilia, Dennis Ulmer, Haau-Sing Li, Raquel Fernández, Barbara Plank, Rico Sennrich, Chrysoula Zerva, and Wilker Aziz. 2023. Uncertainty in natural language generation: From theory to applications. *arXiv preprint arXiv:2307.15703*.

- Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 Conference on Machine Translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.
- Julia Belikova, Evegeniy Beliakin, and Vasily Konovalov. 2024. JellyBell at TextGraphs-17 shared task: Fusing large language models with external knowledge for enhanced question answering. In *Proceedings of TextGraphs-17: Graph-based Methods for Natural Language Processing*, pages 154–160, Bangkok, Thailand. Association for Computational Linguistics.
- Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and Chris Kedzie. 2024. Do androids know they're only dreaming of electric sheep? In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4401–4420, Bangkok, Thailand. Association for Computational Linguistics.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024. INSIDE: LLMs' internal states retain the power of hallucination detection. In *The Twelfth International Conference on Learning Representations, ICLR* 2024, *Vienna, Austria, May* 7-11, 2024.
- Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li, and Yanghua Xiao. 2023. Hallucination detection: Robustly discerning reliable answers in large language models. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 245–255.
- 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.
- Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ranjay Krishna, Yoon Kim, and James R. Glass. 2024. Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps. In *Proceedings of the 2024 Con*ference on Empirical Methods in Natural Language Processing, pages 1419–1436, Miami, Florida, USA. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.

- 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, Bangkok, Thailand. Association for Computational Linguistics.
- Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. Fact-checking the output of large language models via token-level uncertainty quantification. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9367–9385, Bangkok, Thailand. Association for Computational Linguistics.
- Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2023. LM-Polygraph: Uncertainty estimation for language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 446–461, Singapore. Association for Computational Linguistics.
- Javier Ferrando and Elena Voita. 2024. Information flow routes: Automatically interpreting language models at scale. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17432–17445, Miami, Florida, USA. Association for Computational Linguistics.
- Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA. PMLR.
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppl, Preslav Nakov, and Iryna Gurevych. 2024. A survey of confidence estimation and calibration in large language models. 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 6577–6595, Mexico City, Mexico. Association for Computational Linguistics.

- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, et al. 2024. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Jianfeng He, Linlin Yu, Shuo Lei, Chang-Tien Lu, and Feng Chen. 2024a. Uncertainty estimation on sequential labeling via uncertainty transmission. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2823–2835, Mexico City, Mexico. Association for Computational Linguistics.
- Jianfeng He, Xuchao Zhang, Shuo Lei, Zhiqian Chen, Fanglan Chen, Abdulaziz Alhamadani, Bei Xiao, and ChangTien Lu. 2020. Towards more accurate uncertainty estimation in text classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8362–8372, Online. Association for Computational Linguistics.
- Jinwen He, Yujia Gong, Zijin Lin, Cheng'an Wei, Yue Zhao, and Kai Chen. 2024b. LLM factoscope: Uncovering LLMs' factual discernment through measuring inner states. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10218–10230, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.
- Mykyta Ielanskyi, Kajetan Schweighofer, Lukas Aichberger, and Sepp Hochreiter. 2025. Addressing pitfalls in the evaluation of uncertainty estimation methods for natural language generation. In *ICLR Workshop: Quantify Uncertainty and Hallucination in Foundation Models: The Next Frontier in Reliable*
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):1–38.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly

- supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. 2020. Captum: A unified and generic model interpretability library for PyTorch. *arXiv preprint arXiv:2009.07896*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5*, 2023.
- Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym Korablyov, and Yoshua Bengio. 2023. DEUP: Direct epistemic uncertainty prediction. *Transactions on Machine Learning Research*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. Generating with confidence: Uncertainty quantification for black-box large language models. *Transactions on Machine Learning Research*.
- Yu Lu, Jiali Zeng, Jiajun Zhang, Shuangzhi Wu, and Mu Li. 2022. Learning confidence for transformer-based neural machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2353–2364, Dublin, Ireland. Association for Computational Linguistics.
- Andrey Malinin and Mark J. F. Gales. 2021. Uncertainty estimation in autoregressive structured prediction. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. 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*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Viktor Moskvoretskii, Maria Marina, Mikhail Salnikov, Nikolay Ivanov, Sergey Pletenev, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Irina

- Nikishina, and Alexander Panchenko. 2025. Adaptive retrieval without self-knowledge? Bringing uncertainty back home. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 August 1, 2025*, pages 6355–6384. Association for Computational Linguistics.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Finegrained uncertainty quantification for llms from semantic similarities. *Advances in Neural Information Processing Systems*, 37:8901–8929.
- Yookoon Park and David Blei. 2024. Density uncertainty layers for reliable uncertainty estimation. In *International Conference on Artificial Intelligence and Statistics*, pages 163–171. PMLR.
- Xin Qiu and Risto Miikkulainen. 2024. Semantic density: Uncertainty quantification for large language models through confidence measurement in semantic space. In *Advances in Neural Information Processing Systems*, volume 37, pages 134507–134533. Curran Associates, Inc.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J Liu. 2023. Out-of-distribution detection and selective generation for conditional language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.*
- Morgane Rivière, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, et al. 2024. Gemma 2: Improving open language models at a practical size. arXiv preprint arXiv:2408.00118.

- Andrea Santilli, Adam Golinski, Michael Kirchhof, Federico Danieli, Arno Blaas, Miao Xiong, Luca Zappella, and Sinead Williamson. 2025. Revisiting uncertainty quantification evaluation in language models: Spurious interactions with response length bias results. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics* (Volume 2: Short Papers), pages 743–759, Vienna, Austria. Association for Computational Linguistics.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointergenerator networks. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1073– 1083, Vancouver, Canada. Association for Computational Linguistics.
- Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024. Unsupervised real-time hallucination detection based on the internal states of large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14379–14391, Bangkok, Thailand. Association for Computational Linguistics.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328. PMLR.
- Igor Tufanov, Karen Hambardzumyan, Javier Ferrando, and Elena Voita. 2024. LM transparency tool: Interactive tool for analyzing transformer language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 51–60, Bangkok, Thailand. Association for Computational Linguistics.
- Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Daniil Vasilev, Akim Tsvigun, Sergey Petrakov, Rui Xing, Abdelrahman Sadallah, Kirill Grishchenkov, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, Maxim Panov, and Artem Shelmanov. 2025. Benchmarking uncertainty quantification methods for large language models with LM-Polygraph. *Transactions of the Association for Computational Linguistics*, 13:220–248.
- Artem Vazhentsev, Gleb Kuzmin, Artem Shelmanov, Akim Tsvigun, Evgenii Tsymbalov, Kirill Fedyanin, Maxim Panov, Alexander Panchenko, Gleb Gusev, Mikhail Burtsev, Manvel Avetisian, and Leonid Zhukov. 2022. Uncertainty estimation of transformer predictions for misclassification detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8237–8252, Dublin, Ireland. Association for Computational Linguistics.
- Artem Vazhentsev, Gleb Kuzmin, Akim Tsvigun, Alexander Panchenko, Maxim Panov, Mikhail Burt-

- sev, and Artem Shelmanov. 2023. Hybrid uncertainty quantification for selective text classification in ambiguous tasks. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11659–11681, Toronto, Canada. Association for Computational Linguistics.
- Artem Vazhentsev, Lyudmila Rvanova, Ivan Lazichny, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2025. Token-level density-based uncertainty quantification methods for eliciting truthfulness of large language models. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 2246–2262, Albuquerque, New Mexico. Association for Computational Linguistics.
- Yuxia Wang, Daniel Beck, Timothy Baldwin, and Karin Verspoor. 2022. Uncertainty estimation and reduction of pre-trained models for text regression. *Transactions of the Association for Computational Linguistics*, 10:680–696.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. In *Proceedings of the 3rd Workshop on Noisy Usergenerated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.
- Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. The art of abstention: Selective prediction and error regularization for natural language processing. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1040–1051, Online. Association for Computational Linguistics.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, et al. 2024. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115.
- Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. AlignScore: Evaluating factual consistency with a unified alignment function. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.
- Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024. LUQ: Long-text uncertainty quantification for LLMs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5244–5262, Miami, Florida, USA. Association for Computational Linguistics.

Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. 2023. Enhancing uncertainty-based hallucination detection with stronger focus. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 915–932, Singapore. Association for Computational Linguistics.

Xuchao Zhang, Fanglan Chen, Chang-Tien Lu, and Naren Ramakrishnan. 2019. Mitigating uncertainty in document classification. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3126–3136, Minneapolis, Minnesota. Association for Computational Linguistics.

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–38.

A Additional Experimental Results

A.1 Comparison with other UQ Methods

Here, we present the main results for Gemma and Qwen.

TIO M. d 1	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k	Mean	Mean
UQ Method	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.	Acc.	PRR	Rank
MSP	.321	.125	.348	.484	.004	.126	.310	.501	.649	.599	.310	.343	10.27
Perplexity	.325	.263	.078	.449	.397	.278	.314	.234	.660	.578	.256	.348	10.91
Mean Token Entropy	.312	.278	.038	.453	.422	.310	.304	.234	<u>.675</u>	.586	.304	.356	10.45
CCP	.399	.108	.394	.369	.028	.040	.277	.385	.633	.550	.339	.320	12.64
Simple Focus	.248	.174	.277	.521	.170	.270	.335	.523	.656	.570	.280	.366	9.45
Focus	.266	.323	.178	.465	.514	.306	.289	.434	.619	.563	.265	.384	10.18
Lexical Similarity Rouge-L	.059	.186	019	.404	035	.137	.319	.395	.585	.418	.346	.254	14.82
EigenScore	.024	.079	.004	.249	024	.106	.270	.359	.519	.371	.241	.200	19.18
EVL NLI Score entail.	.058	.261	.048	.302	.176	.236	.304	.389	.615	.398	.284	.279	14.91
Ecc. NLI Score entail.	.004	.069	.001	.343	.037	.230	.299	.419	.569	.399	.228	.236	17.45
DegMat NLI Score entail.	.058	.263	.056	.312	.167	.206	.312	.422	.619	.401	.293	.283	13.91
Semantic Entropy	.011	.114	.070	.401	.083	.007	.265	.355	.551	.427	.328	.237	16.64
SAR	.086	.215	.043	.455	.203	.260	.323	.362	.626	.493	.355	.311	11.73
LUQ	.197	.207	.129	.276	.222	.352	.301	.342	.618	.440	.237	.302	13.82
Semantic Density	.172	.198	.124	.313	.272	.448	.401	.463	.654	.295	.183	.320	11.64
Factoscope	.234	113	.036	.021	.471	.077	013	.306	.252	.360	.028	.151	19.00
SAPLMA	.395	.406	.068	.631	.667	.449	015	.455	.498	.422	.447	.402	9.36
MIND	.428	.274	.242	<u>.641</u>	<u>.671</u>	<u>.490</u>	.252	.454	.461	.602	.592	.464	7.55
Sheeps	.523	.327	.315	.640	.436	.409	.270	.455	.511	.491	.654	.457	8.09
LookBackLens	<u>.593</u>	.501	.367	.660	.705	.431	.317	.574	.555	.567	.562	.530	4.64
SATRMD	.506	.330	.347	.525	.612	.315	.359	.595	.681	<u>.615</u>	.443	.484	4.00
TAD	.617	<u>.498</u>	<u>.383</u>	.623	.590	.535	<u>.395</u>	<u>.581</u>	.672	.625	<u>.610</u>	.557	2.36

Table 5: PRR↑ for Gemma-2 9b model for various tasks for the considered sequence-level methods. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

UQ Method	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k	Mean	Mean
	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.	Acc.	PRR	Rank
MSP	.077	.012	.339	.451	.030	088	.291	.551	.610	.654	.268	.291	12.91
Perplexity	.237	.250	.172	.466	.131	.274	.270	.385	.601	.400	.456	.331	10.45
Mean Token Entropy	.233	.280	.149	.475	.143	.356	.263	.342	.603	.225	.469	.322	10.55
CCP	.240	.025	.365	.388	.015	104	.215	.468	.596	.412	.281	.264	14.36
Simple Focus	.109	.116	.191	.496	.021	.093	.321	.536	.620	.550	.310	.306	11.55
Focus	.209	.144	.110	.452	.123	.189	.249	.462	.568	.037	.273	.256	14.55
Lexical Similarity Rouge-L	.122	.057	.122	.370	.075	.159	.297	.507	.531	.274	.511	.275	14.00
EigenScore	.077	010	.073	.374	.018	018	.281	.510	.500	.243	.537	.235	16.18
EVL NLI Score entail.	.139	.145	.068	.294	.122	.306	.329	.519	.571	.236	.372	.282	13.09
Ecc. NLI Score entail.	047	.032	015	.368	.107	.146	.294	.535	.543	.237	.386	.235	15.45
DegMat NLI Score entail.	.138	.145	.075	.332	.122	.300	.329	.540	.574	.235	.402	.290	12.36
Semantic Entropy	.016	.074	.106	.366	.073	.087	.265	.491	.536	.165	.380	.233	17.18
SAR	.128	.129	.107	.445	.088	.185	.318	.526	.585	.288	.459	.296	12.09
LUQ	.228	.170	.131	.265	.096	.322	.337	.449	.580	.321	.331	.294	12.18
Semantic Density	.080	.122	.213	.358	.095	.300	<u>.386</u>	.514	.603	.203	.381	.296	12.27
Factoscope	.185	032	.001	.069	.447	.137	.122	.345	.406	.844	101	.220	17.27
SAPLMA	.245	.326	.009	.345	.018	.321	.001	.374	.497	.418	.440	.272	14.09
MIND	.220	.133	.263	.365	<u>.517</u>	.314	.346	.496	.608	.883	.738	.444	7.36
Sheeps	.361	.313	.258	.487	.391	.476	.357	.487	<u>.663</u>	.883	<u>.710</u>	.490	4.73
LookBackLens	<u>.436</u>	.386	.369	.539	.497	<u>.485</u>	.352	.600	.585	.873	.627	.523	3.55
SATRMD	.338	.322	.254	.525	.362	.254	.315	.547	.623	<u>.885</u>	.566	.454	5.55
TAD	.460	.416	.450	.553	.583	.500	.407	<u>.563</u>	.665	.893	.701	.563	1.27

Table 6: PRR↑ for Qwen-2.5 7b model for various tasks for the considered sequence-level methods. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

A.2 Results Using the ROC-AUC Metric

The results with the ROC-AUC metric are presented in Tables 7 to 9. We obtain discrete versions of the generation quality metrics by thresholding the original continuous values. The thresholds were empirically determined as 0.3 for XSum, SamSum, and CNN/DailyMail; 0.5 for MedQUAD, TruthfulQA, CoQA, SciQ, and TriviaQA; and 0.85 for WMT19. The results align with the trends observed in the PRR metric. Overall, TAD outperforms the second-best method (Sheeps) by 2.4% for LLaMa-3.1 8B, and LookBackLens by 0.1% for Gemma-2 9B, and 2% for Qwen-2.5 7B on average across all datasets.

UQ Method	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	MedQUAD AlignScore	TruthfulQA AlignScore	CoQA AlignScore	SciQ AlignScore	TriviaQA AlignScore	MMLU Acc.	GSM8k Acc.	Mean ROC-AUC	Mean Rank
MSP	.705	.547	.669	.726	.717	.612	.655	.776	.809	.771	.672	.696	11.27
Perplexity	.750	.757	.573	.720	.717	.641	.665	.678	.804	.741	.652	.700	11.50
Mean Token Entropy	.743	.765	.547	.728	.720	.648	.662	.669	.815	.619	.664	.689	12.18
CCP	.727	.525	.676	.671	.720	.556	.633	.704	.824	.709	.678	.674	13.00
Simple Focus	.683	.629	.618	.738	.711	.613	.671	.804	.821	.758	.669	.701	9.82
Focus	.710	.665	.548	.708	.731	.641	.624	.785	.793	.642	.668	.683	12.14
Lexical Similarity Rouge-L	.539	.531	.547	.703	.530	.525	.675	.755	.796	.630	.682	.628	16.00
EigenScore	.521	.548	.552	.668	.543	.537	.658	.710	.775	.603	.677	.617	17.64
EVL NLI Score entail.	.508	.558	.581	.630	.565	.605	.691	.744	.828	.630	.627	.633	14.45
Ecc. NLI Score entail.	.470	.524	.505	.654	.565	.559	.677	.740	.809	.625	.655	.617	17.18
DegMat NLI Score entail.	.509	.564	.583	.640	.572	.606	.689	.778	.834	.633	.637	.640	13.00
Semantic Entropy	.501	.507	.543	.692	.605	.538	.653	.724	.792	.604	.703	.624	17.55
SAR	.548	.580	.548	.727	.641	.564	.683	.796	.825	.661	.689	.660	11.45
LUQ	.587	.656	.591	.629	.554	.580	.684	.774	.823	.637	.605	.647	13.82
Semantic Density	.586	.565	.553	.638	.679	.612	.720	.785	.829	.622	.614	.655	12.73
Factoscope	.707	.575	.494	.592	.832	.541	.513	.698	.705	.820	.558	.640	16.45
SAPLMA	.698	.704	.569	.792	.760	.603	.509	.741	.728	.733	.713	.686	11.91
MIND	.786	.719	.622	.748	.868	.696	.654	.813	.785	.884	.795	.761	6.73
Sheeps	.824	.766	.721	.778	.819	.687	.746	.827	.827	.881	.816	.790	3.09
LookBackLens	.835	.760	.681	.820	.846	.718	.701	.826	.778	.874	.780	.784	4.82
SATRMD	.807	<u>.771</u>	.634	.729	.767	.694	.702	.861	.819	.879	.787	.768	4.73
TAD	.851	.801	.758	.812	.896	.726	<u>.730</u>	.820	.859	.888	<u>.809</u>	.814	1.55

Table 7: ROC-AUC↑ of UQ methods for the Llama-3.1 8b model. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

UQ Method	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	MedQUAD AlignScore	TruthfulQA AlignScore	CoQA AlignScore	SciQ AlignScore	TriviaQA AlignScore	MMLU Acc.	GSM8k Acc.	Mean ROC-AUC	Mean Rank
MSP	.712	.563	.681	.732	.798	.550	.698	.786	.863	.846	.681	.719	9.82
Perplexity	.730	.703	.542	.733	.837	.616	.703	.664	.867	.840	.634	.715	10.73
Mean Token Entropy	.725	.713	.523	.734	.839	.622	.689	.658	.874	.841	.657	.716	10.91
CCP	.748	.535	.697	.679	.810	.520	.669	.711	.857	.816	.699	.704	12.82
Simple Focus	.666	.626	.633	.761	.805	.588	.712	.806	.865	.838	.671	.725	9.36
Focus	.678	.720	.604	.721	.852	.619	.676	.774	.842	.830	.664	.726	10.27
Lexical Similarity Rouge-L	.543	.607	.508	.713	.495	.557	.693	.763	.830	.744	.722	.652	15.27
EigenScore	.520	.529	.491	.639	.478	.550	.673	.729	.797	.725	.662	.618	19.36
EVL NLI Score entail.	.544	.685	.534	.668	.536	.614	.697	.756	.844	.735	.685	.663	14.36
Ecc. NLI Score entail.	.512	.533	.518	.681	.528	.579	.694	.755	.818	.739	.673	.639	16.55
DegMat NLI Score entail.	.543	.688	.539	.671	.544	.608	.694	.771	.845	.724	.691	.665	14.09
Semantic Entropy	.506	.562	.551	.706	.588	.496	.679	.736	.820	.768	.702	.647	15.73
SAR	.574	.638	.532	.745	.641	.615	.700	.771	.851	.795	.726	.690	11.36
LUQ	.631	.647	.571	.648	.521	.663	.684	.755	.844	.745	.661	.670	14.45
Semantic Density	.587	.707	.561	.661	.655	.680	<u>.734</u>	.772	.858	.697	.634	.686	11.91
Factoscope	.677	.409	.543	.529	.925	.537	.493	.715	.640	.718	.523	.610	17.91
SAPLMA	.755	.763	.527	.831	.868	.725	.499	.772	.776	.738	.760	.729	10.00
MIND	.765	.704	.620	.826	.928	.729	.660	.766	.756	.847	.821	.766	7.73
Sheeps	.776	.704	.663	.831	.889	.703	.668	.764	.782	.806	.853	.767	8.45
LookBackLens	.847	.821	.698	.845	.915	.716	.697	.852	.816	.828	.817	.805	4.73
SATRMD	.802	.756	.684	.774	.920	.679	.718	<u>.845</u>	.869	.846	.755	.786	4.09
TAD	.859	<u>.814</u>	.696	.824	.788	.744	.744	.831	<u>.873</u>	.855	.844	.806	3.09

Table 8: ROC-AUC↑ for Gemma-2 9b model for various tasks for the considered sequence-level methods. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

HOM 4. 1	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k	Mean	Mean
UQ Method	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.	Acc.	ROC-AUC	Rank
MSP	.524	.443	.664	.733	.670	.463	.678	.797	.819	.813	.670	.661	13.55
Perplexity	.653	.612	.622	.755	.684	.617	.674	.735	.822	.697	.746	.693	10.27
Mean Token Entropy	.650	.628	.624	.760	.679	.652	.668	.711	.832	.590	.756	.686	10.73
CCP	.650	.445	.689	.688	.665	.452	.639	.751	.821	.673	.679	.650	14.18
Simple Focus	.571	.564	.607	.751	.604	.565	.698	.798	.826	.772	.693	.677	11.36
Focus	.644	.571	.585	.725	.626	.597	.641	.746	.799	.439	.676	.641	15.18
Lexical Similarity Rouge-L	.579	.499	.590	.709	.519	.583	.683	.797	.799	.654	.795	.655	13.73
EigenScore	.541	.532	.570	.702	.526	.502	.670	.795	.780	.640	.812	.643	15.55
EVL NLI Score entail.	.611	.571	.545	.665	.497	.648	.703	.802	.814	.627	.733	.656	13.36
Ecc. NLI Score entail.	.479	.503	.474	.688	.554	.565	.685	.806	.797	.637	.730	.629	16.00
DegMat NLI Score entail.	.611	.575	.547	.675	.506	.642	.698	.807	.819	.631	.743	.659	12.55
Semantic Entropy	.491	.497	.537	.703	.645	.529	.673	.791	.795	.600	.731	.636	16.82
SAR	.589	.546	.596	.743	.597	.621	.695	.808	.821	.663	.773	.677	11.00
LUQ	.668	.535	.587	.647	.479	.657	.699	.776	.814	.673	.721	.660	13.27
Semantic Density	.530	.646	.611	.677	.602	.634	<u>.731</u>	.798	.828	.621	.739	.674	11.09
Factoscope	.617	.472	.506	.540	.735	.546	.585	.706	.716	.909	.409	.613	17.18
SAPLMA	.666	.682	.534	.667	.478	.680	.501	.720	.749	.709	.761	.650	14.18
MIND	.651	.561	.658	.682	.714	.680	.722	.795	.812	.939	.882	.736	8.00
Sheeps	.724	.652	.677	.760	.707	.736	.695	.776	.846	.945	.885	.764	5.00
LookBackLens	<u>.772</u>	.764	<u>.714</u>	.772	.750	.739	.706	.843	.827	.928	.839	<u>.787</u>	2.82
SATRMD	.700	.721	.627	<u>.774</u>	.714	.622	.697	.797	.836	.939	.819	.749	5.73
TAD	.774	<u>.752</u>	.739	.795	.829	<u>.738</u>	.747	<u>.820</u>	.854	.947	.881	.807	1.45

Table 9: ROC-AUC† for Qwen-2.5 7b model for various tasks for the considered sequence-level methods. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

A.3 Generalization to Out-of-Domain Tasks

In this experiment, we examine how our approach can be generalized on the unseen datasets. For each target dataset, we construct a general QA training dataset by sampling 300 instances from the training datasets from each of other QA datasets. Thus, we evaluate TAD that is not trained on the target dataset. We conduct experiments on one dataset from each task: XSum, SamSum, CNN, WMT19, CoQA, SciQ, TriviaQA, MMLU, and GSM8k. We compare the results with three baseline methods: SAR, Semantic Density, and MSP.

Table 3 presents the performance of the supervised methods against the MSP baseline on QA tasks, while Table 10 presents the results when trained on QA datasets and evaluated on summarization and translation tasks. The results demonstrate that TAD consistently outperforms baselines on unseen QA domains, while its generalization across diverse task types remains limited.

UQ Method	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	Mean PRR
MSP	.303	.107	.329	.459	.299
SAR	.052	.166	.049	.435	.176
Semantic Density	.163	.122	.100	.295	.170
Factoscope	.110	.051	072	.083	.043
SAPLMA	050	.052	036	029	016
MIND	086	<u>.185</u>	.064	.158	.080
Sheeps	.111	.098	067	.013	.039
LookBackLens	.165	.201	.005	018	.088
SATRMD	.352	.096	.482	.364	.323
TAD	.259	.184	103	.087	.107

Table 10: PRR↑ for Llama-3.1 8b model for summarization and translation tasks for the considered supervised sequence-level methods trained on the general QA dataset. Unsupervised methods are not included as their performance is not dependent of the training data. Warmer colors indicate better results. The best method is in **bold**, and the second best one is underlined.

A.4 Replacing Attention Weights with Layer Integrated Gradients (LIG) Features in TAD

In this part, we expand our experiments by incorporating the use of Layer Integrated Gradients (LIG: Sundararajan et al. (2017)) as an alternative or addition to attention weights in the TAD method. The LIG features were computed using Captum's (Kokhlikyan et al., 2020) attribute method, where for each predicted token y_i , attributions were calculated with respect to the input and previously generated tokens. Attribution vectors were aggregated across all layers and aligned to match the shape of the attention matrices.

The motivation behind this experiment was to assess whether attribution-based interpretability features, such as LIG, which estimate token importance with respect to model outputs, could serve as a more semantically grounded alternative to raw attention weights. Given the increasing critique of attention as explanation, it was natural to test whether LIG-based representations improve uncertainty modeling.

Table 11 compares the original TAD method with two modified variants: *TAD (LIG)*, which replaces attention weights entirely with LIG attributions, and *TAD (MIX)*, which concatenates LIG attributions with the original attention weights. The results demonstrate that the *TAD (LIG)* method performs the worst across all tasks, particularly on TruthfulQA and SamSum, where it achieves notably low PRR scores. While *TAD (MIX)* significantly outperforms the LIG-only variant, the original TAD method remains superior, achieving the highest average performance across all datasets.

The experiment demonstrates that LIG attributions, while interpretable and semantically grounded, are ineffective as a replacement for attention weights for uncertainty quantification. Furthermore, combining attention weights with LIG attributions can worsen the performance of the TAD method.

A.5 Ablation Studies

Here, we present ablation studies for various numbers of the preceding tokens, different features, and the impact of various layers for the TAD method.

UQ Method	SamSum	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU
e Q Methou	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.
TAD (LIG)	0.246	.252	0.447	0.553	0.669	0.729
TAD (MIX)	0.392	<u>.521</u>	0.510	0.633	0.716	0.789
TAD	0.431	.565	0.509	0.644	0.737	0.806

Table 11: PRR[↑] for Llama-3.1 8b model for various modifications of the TAD method using the LIG features. The best method is in **bold**, the second best is <u>underlined</u>.

UQ Method	Aggregation	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	MedQUAD AlignScore	TruthfulQA AlignScore	CoQA AlignScore	SciQ AlignScore	TriviaQA AlignScore	MMLU Acc.	GSM8k Acc.	Mean PRR
TAD (LinReg)	$\frac{1}{K}\sum_{k=1}^{K} p_k$.550	.535	.444	.592	.624	.463	.392	.488	.632	.696	.557	.543
TAD (LinReg)	$\sum_{k=1}^{K} \log p_k$.438	.208	.422	.358	.444	.311	.287	.474	.604	.724	.355	.420
TAD (MLP)	$\frac{1}{K}\sum_{k=1}^{K} p_k$.538	.529	.445	.578	.526	.460	.388	.477	.634	.717	.537	.530
TAD (MLP)	$\sum_{k=1}^{K} \log p_k$.492	.359	.456	.318	.328	.302	.250	.492	.615	.740	.420	.434

Table 12: Comparison of various considered regression models and aggregation strategies for TAD (PRR↑, Llama-3.1 8b model). Warmer colors indicate better results. The best method is in **bold**, the second best is <u>underlined</u>.

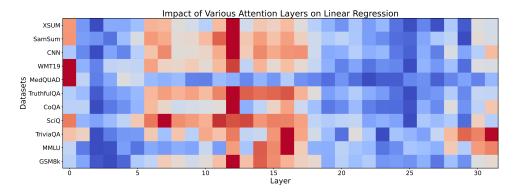


Figure 4: Normalized average weights of linear regression for different attention layers in the TAD method across the considered datasets. Warmer color indicates a higher impact on the TAD performance.

UQ Method	XSum AlignScore	SamSum AlignScore	CNN AlignScore	WMT19 Comet	MedQUAD AlignScore	TruthfulQA AlignScore	GSM8k Acc.	Mean PRR
TAD (1 token)	.554	.538	.437	.563	.496	.407	.513	.519
TAD (2 tokens)	.560	.547	.442	.578	.534	.467	.524	.532
TAD (5 tokens)	.550	.549	<u>.443</u>	<u>.585</u>	<u>.584</u>	.446	<u>.556</u>	.539
TAD (10 tokens)	.550	.535	.444	.592	.624	.463	.557	.545

Table 13: PRR[↑] for Llama-3.1 8b model for various tasks for the various number of preceding tokens for the TAD method. Warmer color indicates better results. The best method is in **bold**, the second best is underlined.

UQ Method	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	GSM8k	Mean
	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	Acc.	PRR
TAD (Sequence-level)	.541	.550	.411	.640	.500	.420	.517	.511
TAD	.550	.535	.444	.592	.624	.463	.557	.538

Table 14: PRR↑ for the modifications of the TAD method for the Llama-3.1 8b model. The best method is in **bold**, the second best is underlined.

UQ Method	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k	Mean
	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.	Acc.	PRR
TAD (probs.)	.554	. 538	.437	.563	.496	.407	.385	.461	.628	.727	.513	.519
TAD (attention)	.549	.532	.453	.590	.624	. 465	.396	.473	.609	.730	.538	.542
TAD (attention+probs.)	.550	.535	.444	.592	.624	.463	.392	.488	.632	.724	.557	.545

Table 15: PRR↑ for Llama-3.1 8b model for various tasks for different features for the TAD method. Warmer color indicates better results. The best method is in **bold**, the second best is <u>underlined</u>.

UQ Method	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k	Mean
	AlignScore	AlignScore	AlignScore	Comet	AlignScore	AlignScore	AlignScore	AlignScore	AlignScore	Acc.	Acc.	PRR
TAD (1 step)	.013	.153	.195	.269	121	.156	.257	.426	.522	.541	.205	.238
TAD (2 step)	.550	. 535	.444	.592	.624	.463	.392	.488	. 632	. 724	.557	.545

Table 16: PRR[↑] for Llama-3.1 8b model for various tasks for the different number of learning steps for the TAD method. Warmer color indicates better results. The best method is in **bold**.

B Computational Resources and Efficiency

All experiments were conducted on a single NVIDIA H100 GPU. On average, training a single model across all datasets took over 750 GPU hours, while inference on the test set took 260 GPU hours.

C Hyperparameters

C.1 Optimal Hyperparameters for TAD

The optimal hyperparameters for TAD for various considered regression models and different aggregation strategies are presented in Tables 17 to 19 for Llama-3.1 8b, Gemma-2 9b, and Qwen-2.5 7b models respectively. These hyperparameters are obtained using cross-validation with five folds using the training dataset. We train a regression model on k-1 folds of the training dataset and estimate uncertainty on the remaining fold. The optimal hyperparameters are selected according to the best average PRR for AlignScore. Finally, we use these hyperparameters to train the regression model on the entire training set.

The hyperparameter grid for the linear regression is the following:

L2 regularization: [1e+1, 1, 1e-1, 1e-2, 1e-3, 1e-4].

The hyperparameter grid for the MLP is the following:

Num. of layers: [2, 4];

Num. of epochs: [10, 20, 30]; Learning rate: [1e-5, 3e-5, 5e-5];

Batch size: [64, 128]. For both models, we include aggregation strategies into the hyperparameter grid

for the final configuration.

UQ Method	Aggregation	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k
TAD (MLP)												4, 30, 1e-05, 0, 128
TAD (MLP)	$\sum_{k=1}^{K} \log p_k$	4, 30, 3e-05, 0, 128	4, 30, 5e-05, 0, 64	4, 30, 5e-05, 0, 64	2, 30, 3e-05, 0, 128	4, 30, 1e-05, 0, 128	4, 30, 5e-05, 0, 128	2, 30, 5e-05, 0, 64	4, 30, 5e-05, 0, 64	4, 30, 3e-05, 0, 128	4, 30, 5e-05, 0, 128	4, 30, 1e-05, 0, 128
TAD (LinReg)	$\frac{1}{K}\sum_{k=1}^{K} p_k$	1	1	10.0	1	0.001	0.1	1	0.01	10.0	1	10.0
TAD (LinReg)	$\sum_{k=1}^{K} \log p_k$	10.0	0.01	1	0.001	0.001	1	1	1	10.0	1	0.1

Table 17: Optimal values of the hyper-parameters for the TAD methods for the Llama-3.1 8b model.

UQ Method XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k
TAD (LinReg) 0.01	1	1	1	0.001	0.1	10.0	0.1	10.0	1	0.1

Table 18: Optimal values of the hyper-parameters for the final configuration of the TAD method for the Gemma-2 9b model. For CNN, SciQ, and MMLU, $\sum_{k=1}^K \log p_k$ is the best aggregation method, whereas $\frac{1}{K} \sum_{k=1}^K p_k$ performs best on all other datasets.

UQ Method	XSum	SamSum	CNN	WMT19	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA	MMLU	GSM8k
TAD (LinReg)	0.01	1	10.0	1	0.001	0.1	10.0	0.1	10.0	1	0.1

Table 19: Optimal values of the hyper-parameters for the final configuration of the TAD method for the Qwen-2.5 7b model. For MMLU, $\sum_{k=1}^{K} \log p_k$ is the best aggregation method, whereas $\frac{1}{K} \sum_{k=1}^{K} p_k$ performs best on all other datasets.

C.2 LLM Generation Hyperparameters

Dataset	Task	Max Input Length	Generation Length	Temperature	Тор-р	Do Sample	Beams	Repetition Penalty
XSum			56					
SamSum	TS		128					
CNN			128					1.2
WMT19	MT		107					
MedQUAD	QA		128					
TruthfulQA	_	-	128	1.0	1.0	False	1	1
GSM8k	Long answer		256					
CoQA	0.4		20					
SciQ	QA		20					
TriviQA	Short answer		20					
MMLU	MCQA		3					

Table 20: Values of the text generation hyper-parameters for all LLMs used in our experiments.

D Dataset Statistics

Statistics about the datasets are provided in Table 21. For TS, we experiment with CNN/DailyMail (See et al., 2017), XSum (Narayan et al., 2018), and SamSum (Gliwa et al., 2019). For the long answer QA task, we use MedQUAD (Abacha and Demner-Fushman, 2019), which consists of real medical questions, TruthfulQA (Lin et al., 2022), which consists of questions that some people would answer incorrectly due to a false belief or a misconception, and GSM8k (Cobbe et al., 2021) with a grade school math questions. For the QA task with short answers, we follow previous work on UQ (Kuhn et al., 2023; Duan et al., 2024; Lin et al., 2024) and we use three datasets: SciQ (Welbl et al., 2017), CoQA (Reddy et al., 2019), and TriviaQA (Joshi et al., 2017). For multiple-choice QA, we use MMLU (Hendrycks et al., 2021), a widely used benchmark for evaluating LLMs. For MT, we use WMT19 (Barrault et al., 2019), focusing on translations from German to English.

Task	Dataset	N-shot	Train texts for TAD	Evaluation texts
Text	CNN/DailyMail	0	500	1,000
Summarization	XSum	0	1,000	2,000
Summarization	SamSum	0	2,000	819
MT	WMT19 De-En	0	2,000	2,000
OA	MedQUAD	5	500	1,000
	TruthfulQA	5	408	409
Long answer	GSM8k	5	700	1,319
-	SciQ	0	2,000	1,000
QA Short answer	CoQA	all preceding questions	2,000	2,000
	TriviaQA	5	2,000	2,000
MCQA	MMLU	5	2,000	2,000

Table 21: Statistics about the datasets used for evaluation.

E Generating Training Data for TAD

Algorithm 1: Generating training data for TAD Data: Input prompt \mathbf{x}_k , LLM generation $\mathbf{y}_k = \mathbf{y}_{1:n_k}$, token probabilities $p(y_i \mid \mathbf{y}_{< i}, \mathbf{x}_k)$, number of preceding tokens N, vector of LLM attention weights $a_{i,i-l}$ from the (i-l)-th token to the i-th token from all layers and heads, and step of the training procedure jResult: Feature vectors z_i^k , $k = 1 \dots K$, $i = 2 \dots n_k$ // Estimate unconditional probability for the first token 1 $\hat{p}_k(y_1) = \sin(\mathbf{y}_k, \mathbf{y}_k^*)$;

```
2 for i \leftarrow 2 to n_k do
        // Construct token-level features
        if i == 1 then
3
             // On the first training step, we use only probabilities as features
             z_i^k \leftarrow \bigoplus_{l=1}^{\min\{N,i-1\}} \left[ p(y_{i-l} \mid \mathbf{y}_{< i-l}, \mathbf{x}_k) \right] \oplus \left[ p(y_i \mid \mathbf{y}_{< i}, \mathbf{x}_k) \right];
 4
5
             // On the next training steps, we use all features
            z_i^k \leftarrow \bigoplus_{l=1}^{\min\{N,i-1\}} \left[ p(y_{i-l} \mid \mathbf{y}_{< i-l}, \mathbf{x}_k), \ \hat{p}_k(y_{i-l}), \ a_{i,i-l} \right] \oplus \left[ p(y_i \mid \mathbf{y}_{< i}, \mathbf{x}_k) \right];
        // If N>i-1, we pad z_i^k with zeros to ensure they have the same length
        if i - 1 < N then
         z_i^k \leftarrow z_i^k \oplus \mathbf{0}_{(2+|a_{i.i-l}|)(N-i+1)};
        // Estimate token-level unconditional probability
        if j == 1 then
             // On the first training step, we use ground truth
             \hat{p}_k(y_i) = \sin(\mathbf{y}_k, \mathbf{y}_k^*);
10
        else
11
             // On the next training steps, we use trained function C(\cdot)
          \hat{p}_k(y_i) = C(z_i^k);
```