

Trainable Attention-based Conditional Dependency for Uncertainty Quantification of Large Language Models

Anonymous ACL submission

Abstract

Uncertainty quantification (UQ) is a perspective approach to detecting Large Language Model (LLM) hallucinations and low quality output. In this work, we address one of the challenges of UQ in generation tasks that arises from the conditional dependency between the generation steps of a LLM. We propose to learn this dependency from data. We train a regression model, which target variable is the gap between the conditional and the unconditional generation confidence. During LLM inference, we use this learned conditional dependency model to modulate the uncertainty of the current generation step based on the uncertainty of the previous step. Our experimental evaluation on nine datasets and three LLMs shows that the proposed method is highly effective for uncertainty quantification, achieving substantial improvements over rivaling approaches.

1 Introduction

Uncertainty quantification (UQ) (Gal and Ghahramani, 2016; Baan et al., 2023; Geng et al., 2023; Fadeeva et al., 2023) is of growing interest in the Natural Language Processing (NLP) community for dealing with Large Language Models (LLMs) hallucinations (Fadeeva et al., 2024) and low quality generations (Malinin and Gales, 2021) in an efficient manner. For example, high uncertainty could serve as an indicator that the entire generation should be discarded (selective generation) to prevent potential harms to users, or that part of the generation should be highlighted to the user as untrustworthy.

There are many approaches for detecting hallucinations and low-quality outputs of LLMs (Manakul et al., 2023; Min et al., 2023; Chen et al., 2023). However, the majority of them leverage external knowledge 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 models inherently contain information about their own knowledge limitations, and that there should be an efficient way to access this information, which can enable LLM-based applications that are both safer and easy to use in practice.

For general classification and regression tasks and for text classification in particular, there is a well-developed battery of UQ techniques (Zhang et al., 2019; He et al., 2020; Xin et al., 2021; Wang et al., 2022; Vazhentsev et al., 2023). For text generation tasks, UQ is much more complicated. The complexity is multifold: (1) there is an infinite number of possible generations, which complicates the normalization of the uncertainty scores; (2) in the general case, there are an infinite number of correct answers; (3) decisions are generally based on imprecise sampling and inference algorithms such as beam search; (4) there is not one, but multiple predictions, and the uncertainty of these predictions need to be aggregated; and (5) finally, the predictions at each generation step are not conditionally independent (Zhang et al., 2023).

This last problem is the focus of the present work. During generation, the LLM conditions on the previously-generated tokens. Thus, if the LLM has hallucinated and generated an incorrect claim at the beginning of the sequence, all subsequently generated claims might also be incorrect. Even in the case when the first claim was generated with high uncertainty, this is not taken into account during the subsequent generation process. This means that while the first error could be implicitly recognized as such with high uncertainty, all subsequent mistakes will be overlooked, because the generation process conditioned on this error will be very confident.

Below, we suggest a theoretically-motivated data-driven solution to this problem. We note that the attention between generated tokens pro-

041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

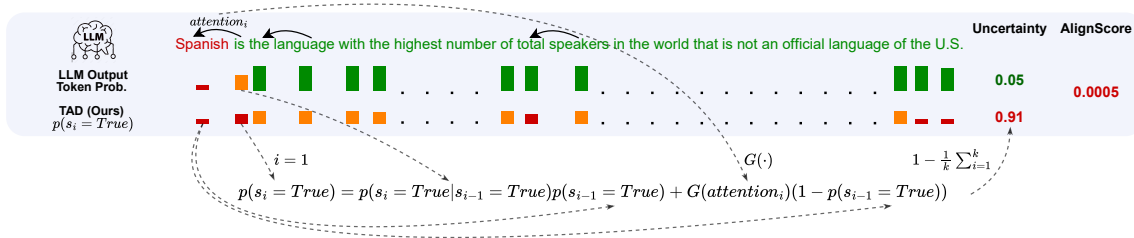


Figure 1: An illustration of the proposed method TAD. The figure depicts generated tokens, uncertainty scores, and probabilities assigned by a LLM (represented with bars). The output was generated by Gemma 7b for the question *What is the language with the highest number of total speakers in the world that is not an official language of the U.S.?* The LLM starts with generating a 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 $G(\cdot)$ based on attention, resulting in a high overall uncertainty for the generated answer.

vides 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 particular algorithmic function would be too difficult to engineer, and thus we propose to learn this dependency from data. For this purpose, we generated a training dataset with a target variable that represents the gap between the conditional generation confidence and the unconditional generation confidence. Using attention-based features, we trained an ML-based regression model to predict this gap that is further used for modifying the certainty of the current generation. We use attention-based features to ensure the generalizability of such an approach, supporting the training of a robust conditional dependency model. We call the proposed approach *trainable attention-based dependency (TAD)*. Figure 1 illustrates the idea behind the proposed method on the real output of an LLM. Our extensive experiments demonstrate that TAD offers substantial improvements in UQ over the baselines in tasks where the LLM is required to generate long sequences.

The contributions of this work are as follows:

- A new data-driven approach to uncertainty quantification that models the conditional dependency between the individual token predictions of an LLM.
- A computationally-efficient implementation of the method that leverages simple linear regression, making it practical for real-world applications based on LLMs.
- An empirical demonstration that our proposed method outperforms previous approaches across nine datasets and three LLMs.

2 Related Work

With the advent of LLMs, UQ has become an urgent research problem in NLP. As previously mentioned, this area not only offers promising practical benefits, but it also presents several intriguing research challenges. The majority of the work on UQ has been unsupervised, with a smaller number of recently-proposed supervised methods.

Unsupervised UQ methods. The challenge of multiple predictions, which we mentioned above, has been previously tackled by aggregating the logits of the generated tokens in various ways and by adapting information-based UQ techniques. Fomicheva et al. (2020) experimented with perplexity and mean token entropy for MT quality estimation. Takayama and Arase (2019) adapted point-wise mutual information (PMI), and van der Poel et al. (2022) extended this approach to conditional PMI. The advantages of these techniques are their simplicity, usually minimal computational overhead, and robust performance. A well-known approach to UQ in general is ensembling (Lakshminarayanan et al., 2017) and Monte Carlo (MC) dropout (Gal and Ghahramani, 2016). Malinin and Gales (2021) and Fomicheva et al. (2020) adapted it to sequence generation problems. In particular, 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).

The problem of multiple correct generations was explicitly addressed in (Kuhn et al., 2023; Nikitin et al., 2024; Cheng and Vlachos, 2024) and in a series of black-box generation methods (Lin et al., 2023). The main idea is to sample multiple genera-

tions from a LLM, extract semantically equivalent clusters from the outputs, and to analyze the diversity of the generated meanings instead of the surface forms.

Fadeeva et al. (2024) addressed the problem of multiple sources of uncertainty present in the LLM probability distribution that are irrelevant for hallucination detection and fact-checking. In addition to dealing with multiple correct generations, they also suggested mitigating the influence of the uncertainty related to the type of generated claims.

Zhang et al. (2023) and Duan et al. (2023) emphasized that not all generated tokens should contribute to the uncertainty score for the entire generated text, and proposed various heuristics to select only relevant tokens.

Zhang et al. (2023) also modeled the conditional dependency between the generation steps by adding a penalty to an uncertainty score that depends on the uncertainties of previously-generated tokens in the form of max-pooled attention to corresponding tokens from the current step.

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 we propose a data-driven approach to 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 the computer vision domain (Lahlou et al., 2022; Park and Blei, 2024). One of their key benefits is computational efficiency.

A handful of papers applied this approach to text generation tasks. Lu et al. (2022) proposed to train a regression head of a model 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 trained an additional head by modifying the loss function and adding a regularizer. However, their approach is only applicable when fine-tuning language models for Machine Translation (MT), and is not suitable for general-purpose instruction-tuned 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. For this purpose, they annotated a dataset of true and false statements, and used forced LM decoding to generate them. They evaluated the ability of the trained MLP to classify the statements as true or false, and demonstrated that it outperformed other supervised baselines and few-shot prompting of the LLM itself. However, due to the reliance on forced decoding, their experimental setup is far from real-world hallucination detection, where an LLM can perform unrestricted generation. Another limitation is that their method can provide veracity scores only for the entire generated text.

Unlike these methods, besides learning uncertainty scores directly from data, we also learn the conditional dependency between the generation steps. Our method is also flexible as it can be used on various levels: for the entire text, at the sub-sentence level, or for individual tokens.

3 Trainable Attention-Based Conditional Dependency

In this section, we present our approach to learn the conditional dependency between the generation steps and our UQ method based on it.

3.1 Theoretical Background and Motivation

When an LLM generates a sequence of tokens t_i , it provides us a conditional probability distribution $p(t_i|t_{<i})$. This essentially means the LLM believes that everything generated so far is correct, which might not be the case. In practice, we would like to somehow propagate its uncertainty from previous generation steps.

Assume for simplicity that we already have some statements s_1, s_2, \dots, s_n and a prompt x , and we have trained a generative LLM to predict the probability of truthfulness of the statements ('T' or 'F') via a Markov process. At each step the LLM provides us $p(s_1 = T | x), p(s_2 = T | s_1 = T), \dots, p(s_n = T | s_{n-1} = T)$. These probability distributions are conditionally dependent on the previous ones. However, to estimate the correctness of some statement s_i , we need to obtain an *unconditional probability* $p(s_i = T)$. The LLM does not provide us such probability during standard generation process. There are some heuristic techniques such as P(true) (Kadavath et al., 2022) that can estimate the unconditional probability through

rerunning LLM on the generated text. However, it introduces expensive overhead, which approximately doubles the generation time.

We would like to have a computationally efficient approach that does not need rerunning the LLM. Let us expand $p(s_i = T)$ according to the formula of full probability and express it using conditional probability:

$$\begin{aligned}
 p(s_i = T) &= p(s_i = T, s_{i-1} = T) + p(s_i = T, s_{i-1} = F) \\
 &= p(s_i = T | s_{i-1} = T) p(s_{i-1} = T) + \\
 &+ p(s_i = T | s_{i-1} = F) p(s_{i-1} = F) \\
 &= p(s_i = T | s_{i-1} = T) p(s_{i-1} = T) + \\
 &+ p(s_i = T | s_{i-1} = F)(1 - p(s_{i-1} = T)).
 \end{aligned} \tag{1}$$

In the obtained formula, $p(s_i = T | s_{i-1} = T)$ is what the LLM provides during the current generation step. Consider that we know $p(s_{i-1} = T)$ as it is calculated on the previous generation step. We still do not know the remaining term: $p(s_i = T | s_{i-1} = F)$. Let us express it from the equation:

$$\begin{aligned}
 p(s_i = T | s_{i-1} = F) & \tag{2} \\
 = \frac{p(s_i = T) - p(s_i = T | s_{i-1} = T) p(s_{i-1} = T)}{1 - p(s_{i-1} = T)}.
 \end{aligned}$$

This expression still requires $p(s_i = T)$, which is not known during the inference. However, we can replace it with some surrogate and use this expression to approximate $p(s_i = T | s_{i-1} = F)$ with a trainable model $G(Atten_i, p(s_{i-1} = T), p(s_i = T | s_{i-1} = T))$. This function in fact measures the conditional dependency of the current generation step i on the previous one $i - 1$. For model features, we suggest using attention from the step i to $i - 1$: $Atten_i$. The training data for this model could be obtained using Equation (2) in the “offline” mode, where we do not care about efficiency of obtaining $p(s_i = T)$. We also note that if the implementation of G is a linear regression or a small neural network, it will not introduce much overhead to compute during the inference of the main LLM.

Finally, to obtain the confidence estimate, we replace $p(s_i = T | s_{i-1} = F)$ with G in Equation (1):

$$\begin{aligned}
 p(s_i = T) &= p(s_i = T | s_{i-1} = T) p(s_{i-1} = T) \\
 &+ G(Atten_i, p(s_{i-1} = T), p(s_i = T | s_{i-1} = T)) \\
 &\cdot (1 - p(s_{i-1} = T)).
 \end{aligned} \tag{3}$$

We note that, in order to implement G , we need an effective way of obtaining unconditional probabilities $p(s_i = T)$ and we also need to deal with the fact that real LLMs produce actually tokens. We address these problems and suggest the implementation of G below.

3.2 Implementation

Despite some strong assumptions, we argue that the presented motivation could be applied to individual tokens t_1, t_2, \dots, t_n as well. We implement the proposed method for the token-level uncertainty scores and then we aggregate these token-level scores into a score for the whole sequence.

Obtaining unconditional probability. To obtain the surrogate for the unconditional confidence $\hat{p}(t_i)$ for a generated token t_i during the training phase, we use two strategies. The first one relies solely on the strict criterion of the presence of an existing token t_i in the ground truth text y :

$$\hat{p}(t_i) = \begin{cases} 1, & t_i \in y, \\ 0, & \text{otherwise.} \end{cases} \tag{4}$$

The second strategy also leverages AlignScore (Zha et al., 2023) $sim(\cdot, \cdot)$ between the generated text \tilde{y} and the ground-truth y :

$$\hat{p}(t_i) = \begin{cases} \frac{1+sim(\tilde{y}, y)}{2}, & t_i \in y, \\ sim(\tilde{y}, y), & \text{otherwise.} \end{cases} \tag{5}$$

This strategy aims to correct the target when a generated token is not present in the expected text, but the AlignScore is high, indicating that the generated text has the similar meaning as the training sentence. In the inverse situation, when the token is present, but the whole generation according to AlignScore is wrong, it penalizes the target.

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 x_k and the target text y_k , using a LLM, we generate a text \tilde{y}_k of some length n_k and token probabilities $p(t_i | t_{<i})$.
2. For the first generated token t_1 in each text, we define its unconditional confidence as a ground truth surrogate $p(t_1) = \hat{p}(t_1)$ according to formulas (4) or (5).
3. For each generated token $t_i, i = 2, \dots, n_k$:
 - (a) We obtain $p(t_{i-1})$ from the previous generation step.

(b) We define its unconditional confidence as a ground truth surrogate $p(t_i) = \hat{p}(t_i)$ according to equations (4) or (5).

(c) We compute the target variable for the function G using equation (2):

$$\tilde{G}_i = \frac{p(t_i) - p(t_i | t_{<i}) p(t_{i-1})}{1 - p(t_{i-1})}.$$

As a result, for each instance in the training dataset, we generate a sequence of target variables \tilde{G}_i^k $k = 1, \dots, K, i = 1, \dots, n_k$. We further train the model G on these targets.

Model for G and its training procedure. We experiment with several regression models for TAD: linear regression (LinReg), CatBoost regression (Prokhorenkova et al., 2018), and a multi-layer perceptron (MLP). The hyperparameters 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-AlignScore. Finally, we use these values to train the regression model on the full training set. The selected hyperparameters 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. The features are used to compute G and a confidence score based on Equation (3).

4 Experiments

4.1 Experimental Setup

For 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 or sent to a more advanced LLM.

Metrics. Following previous work on UQ in text generation (Malinin and Gales, 2021; Fadeeva et al., 2023), 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 metric. 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 selective generation. We use ROUGE-L, Accuracy, and AlignScore (Zha et al., 2023) as generation quality metrics.

Datasets. We consider three text generation tasks: text summarization (TS), QA with long free-form answers, and QA with free-form short answers, and for each task, we consider three datasets. Statistics about the datasets are provided in Table 18 in Appendix D. For TS, we experiment with CNN/DailyMail (See et al., 2017), XSum (Narayan et al., 2018) (summarization of news articles), and SamSum (Gliwa et al., 2019) (summarization of dialogues). For the long answer QA task, we use PubMedQA (Jin et al., 2019), a QA dataset in the biomedical domain, with the task to answer biomedical research questions using the corresponding abstracts. We further use MedQUAD (Abacha and Demner-Fushman, 2019), which consists of real medical questions, and TruthfulQA (Lin et al., 2022), which consists of questions that some people would answer incorrectly due to a false belief or a misconception. For the QA task with short answers, we follow previous work on UQ (Kuhn et al., 2023; Duan et al., 2023; Lin et al., 2023) and we use three datasets: SciQ (Welbl et al., 2017), CoQA (Reddy et al., 2019), and TriviaQA (Joshi et al., 2017).

LLMs. We experiment with three LLMs: Gemma 7b (Mesnard et al., 2024), LLaMA 8b v3, and StableLM 12b v2 (Bellagente et al., 2024). The inference hyperparameters of the LLMs are given in Table 17 in Appendix C.2.

UQ baselines. We compare TAD to 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 more complex techniques, considered to be state-of-the-art UQ methods for LLMs: Lexical Similarity based on ROUGE-L (Fomicheva et al., 2020), Monte Carlo Sequence Entropy (MC SE), Monte Carlo Normalized Sequence Entropy (MC NSE; Kuhn et al. (2023)), black-box methods (NumSemSets, DegMat, Eccentricity, EigVallaplacian; Lin et al. (2023)), Semantic Entropy (Kuhn et al., 2023), hallucination detection with stronger focus (Focus; Zhang et al. (2023)), and Shifting Attention to Relevance (SAR; Duan et al. (2023)).

UQ Method	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
MSP	<u>-329</u>	<u>-116</u>	.234	.177	<u>-039</u>	<u>.043</u>	<u>-455</u>	<u>-154</u>	<u>-454</u>	<u>.008</u>	<u>.520</u>	<u>.268</u>	<u>.699</u>	<u>.626</u>	<u>.806</u>	<u>.744</u>	<u>.828</u>	<u>.805</u>	8.61
Perplexity	<u>-358</u>	<u>-179</u>	.206	<u>.291</u>	.071	<u>-012</u>	<u>-527</u>	<u>-159</u>	<u>.801</u>	<u>.346</u>	<u>.381</u>	<u>.318</u>	<u>.458</u>	<u>.439</u>	<u>-.321</u>	<u>-.399</u>	<u>.820</u>	<u>.791</u>	7.78
Mean Token Entropy	<u>-350</u>	<u>-181</u>	.172	<u>.281</u>	.082	<u>-017</u>	<u>-524</u>	<u>-147</u>	<u>.776</u>	<u>.330</u>	<u>.228</u>	<u>.290</u>	<u>.327</u>	<u>.339</u>	<u>-.398</u>	<u>-.398</u>	<u>.806</u>	<u>.786</u>	8.94
Focus	<u>-324</u>	<u>-161</u>	.169	<u>.232</u>	.023	<u>.008</u>	<u>-357</u>	<u>-146</u>	<u>-408</u>	<u>-100</u>	<u>.306</u>	<u>.298</u>	<u>.322</u>	<u>.280</u>	<u>-.098</u>	<u>.070</u>	<u>.651</u>	<u>.702</u>	13.00
NumSemSets	.054	.049	.176	.176	.029	.052	.041	.017	<u>-067</u>	<u>.047</u>	<u>.132</u>	<u>.231</u>	<u>.203</u>	<u>.349</u>	<u>.132</u>	<u>.275</u>	<u>.677</u>	<u>.714</u>	10.72
DeqMat	.025	.060	.141	.161	.072	.088	.028	.008	<u>-063</u>	<u>.087</u>	<u>.211</u>	<u>.285</u>	<u>.345</u>	<u>.496</u>	<u>.401</u>	<u>.353</u>	<u>.740</u>	<u>.770</u>	8.61
Eccentricity	<u>-055</u>	<u>.010</u>	<u>.059</u>	<u>.052</u>	<u>.028</u>	<u>-005</u>	<u>-016</u>	<u>-011</u>	<u>-144</u>	<u>.027</u>	<u>-116</u>	<u>-213</u>	<u>.514</u>	<u>.559</u>	<u>.487</u>	<u>.570</u>	<u>.737</u>	<u>.739</u>	11.11
EigVallLaplacian	.024	.063	.140	.156	.071	.087	.016	.004	<u>-155</u>	<u>.064</u>	<u>.200</u>	<u>.279</u>	<u>.479</u>	<u>.538</u>	<u>.507</u>	<u>.603</u>	<u>.727</u>	<u>.760</u>	9.00
Lexical Similarity	.076	<u>-024</u>	.256	.233	.108	.066	.068	.023	<u>.240</u>	<u>-024</u>	<u>.145</u>	<u>.117</u>	<u>.504</u>	<u>.499</u>	<u>.488</u>	<u>.538</u>	<u>.730</u>	<u>.734</u>	8.78
MC NSE	<u>-005</u>	<u>-023</u>	.212	.195	.108	<u>.102</u>	.074	.012	<u>-000</u>	<u>.011</u>	<u>.076</u>	<u>.221</u>	<u>.440</u>	<u>.432</u>	<u>.357</u>	<u>.398</u>	<u>.727</u>	<u>.715</u>	10.00
MC SE	.035	<u>-001</u>	.251	.195	.123	.086	<u>-014</u>	<u>-007</u>	<u>-099</u>	<u>.013</u>	<u>.160</u>	<u>.141</u>	<u>.553</u>	<u>.514</u>	<u>.542</u>	<u>.557</u>	<u>.723</u>	<u>.712</u>	9.11
Semantic Entropy	.034	.001	.250	.195	.110	.082	<u>-019</u>	<u>-003</u>	<u>-097</u>	<u>.019</u>	<u>.158</u>	<u>.159</u>	<u>.583</u>	<u>.566</u>	<u>.589</u>	<u>.605</u>	<u>.752</u>	<u>.745</u>	8.28
SentenceSAR	<u>-077</u>	<u>-037</u>	.168	.133	.061	.090	<u>-072</u>	<u>-033</u>	<u>-221</u>	<u>.013</u>	<u>.305</u>	<u>.199</u>	<u>.643</u>	<u>.665</u>	<u>.700</u>	<u>.692</u>	<u>.792</u>	<u>.786</u>	9.06
SAR	.042	<u>-006</u>	.248	.245	.123	<u>.103</u>	.111	.014	.066	.035	<u>.155</u>	<u>.263</u>	<u>.477</u>	<u>.503</u>	<u>.453</u>	<u>.515</u>	<u>.769</u>	<u>.770</u>	7.11
TAD (LinReg)	<u>.502</u>	<u>.257</u>	<u>.329</u>	<u>.263</u>	<u>.177</u>	<u>.078</u>	<u>.576</u>	<u>.242</u>	<u>.287</u>	<u>.376</u>	<u>.563</u>	<u>.294</u>	<u>.671</u>	<u>.608</u>	<u>.820</u>	<u>.751</u>	<u>.782</u>	<u>.760</u>	3.00
TAD (LinReg+AlignScore)	<u>.541</u>	<u>.380</u>	<u>.353</u>	<u>.349</u>	<u>.146</u>	<u>.092</u>	<u>.007</u>	<u>.064</u>	<u>.491</u>	<u>.472</u>	<u>.505</u>	<u>.368</u>	<u>.621</u>	<u>.600</u>	<u>.834</u>	<u>.777</u>	<u>.784</u>	<u>.766</u>	2.89

Table 1: PRR \uparrow of UQ methods for the Gemma 7b model. Warmer colors indicate better results. The best method is in bold, the second best is underlined.

For these methods, we generate five samples.

4.2 Main Results

Fine-grained comparison with the baselines.

Tables 1, 8 and 9 in Appendix A present the results for Gemma 7b, Llama 8b v3, and StableLM 12b v2 models respectively.

We can see that for all summarization datasets, in the majority of cases, TAD outperforms the state-of-the-art methods by a large margin in terms of both considered metrics. The only exception is the case of PRR-AlignScore for StableLM on the XSum dataset, where SAR and Lexical Similarity are marginally better. At the same time, TAD confidently outperforms them in terms of PRR-ROUGE-L. In experiments with two other models on XSum, TAD also demonstrates large improvements in terms of both metrics over the baselines, which typically perform no better than a random choice. For example, TAD LinReg+AlignScore outperforms the second best baseline by .317 PRR-AlignScore and by .465 PRR-ROUGE-L absolute.

For QA with long answer datasets (PubMedQA, MedQUAD, and TruthfulQA), we see that TAD also confidently outperforms the baselines for all considered settings except for the experiment on TruthfulQA with LLaMA 8b v3 and for PRR-ROUGE-L measured on MedQUAD for Gemma. For example, in the experiment with LLaMA 8b v3 on PubMedQA, TAD outperforms the second best baseline – Perplexity by .190 of PRR-ROUGE-L and by .187 of PRR-AlignScore. For StableLM, the improvement is .049 of PRR-ROUGE-L and .083 of PRR-AlignScore. Additionally, we can see that on this task, the majority of sophisticated UQ baselines consistently fall behind simple techniques.

Finally, for QA with short answers (CoQA, SciQ, and TriviaQA), we can see that TAD notably outperforms baselines for all considered LLMs only on the SciQ dataset. TAD also marginally out-

performs baselines in the experiments on CoQA with StableLM and Llama 8b v3. The lower performance on tasks with short answers is expected, since TAD primarily aims at improving the performance for tasks with long generations and complex conditional dependencies. Moreover, we can see that in the short-answer setting on TriviaQA and CoQA, the simplest baseline MSP demonstrates very strong performance, which is often the best.

When comparing the two strategies for obtaining the unconditional probability during training, we see that adding AlignScore usually helps for summarization tasks, but it could negatively impact the performance for QA.

Overall results. Table 2 presents the mean rank of each method aggregated over all datasets for each model separately. The lower rank is better. The column “Mean Rank” corresponds to 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 that some baseline methods might show good results in several individual cases, they usually are quite unstable resulting in poor overall ranking. At the same time, TAD demonstrates more robust improvements across multiple tasks and LLMs, making it a better choice overall.

Generalization of TAD on unseen datasets. Tables 3, 10 and 11 in Appendix A.2 compare the results of TAD trained on a single in-domain training dataset to the results of TAD trained on all training datasets except one that represents the in-domain dataset for testing (we designate it as Gen TAD). This setting evaluates the out-of-domain performance of TAD. We can see that TAD without the AlignScore target demonstrates good general-

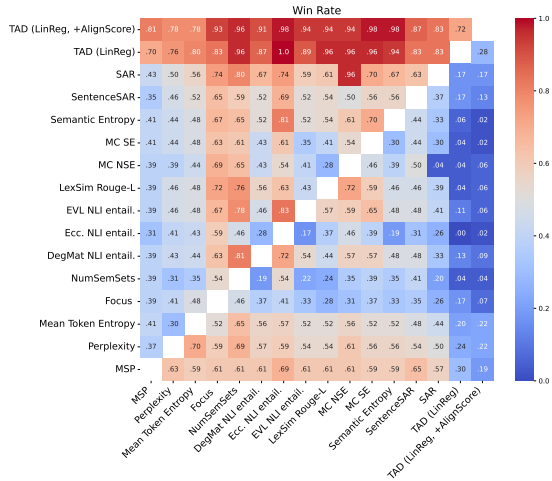


Figure 2: Summary of 54 experimental setups with various models and datasets. Each cell in the diagram presents the fraction of experiments where a method from a row outperforms a method from a column. Warmer colors indicate better results.

UQ Method	Gemma 7b	Llama-3 8b	StableLM 12b	Mean Rank
MSP	8.61	7.17	6.83	4.50
Perplexity	7.78	8.44	8.33	5.33
Mean Token Entropy	8.94	9.11	9.00	9.00
Focus	13.00	9.50	10.50	13.67
NumSemSets	10.72	10.78	12.83	15.00
DegMat	8.61	8.83	9.33	8.17
Eccentricity	11.11	11.33	11.61	15.33
EigValLaplacian	9.00	7.94	8.78	7.67
Lexical Similarity	8.78	9.22	8.56	8.33
MC NSE	10.00	10.72	10.22	13.00
MC SE	9.11	10.22	10.67	13.00
Semantic Entropy	8.28	9.06	9.06	7.67
SentenceSAR	9.06	9.39	8.22	9.00
SAR	7.11	7.78	6.33	3.33
TAD (LinReg)	3.00	3.72	3.50	2.00
TAD (LinReg+AlignSc.)	2.89	2.78	2.22	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.

ization for QA with long answers. Despite the results degrade on the unseen dataset, TAD confidently outperforms other baselines. Adding the AlignScore for QA worsens the results probably due to overfitting.

For the TS task, on the contrary, adding AlignScore helps to achieve some generalization. The results substantially degrade, but are still better than for other baselines. On the short-answer QA task, training on out-of-domain data slightly improves PRR-Accuracy. More details about these experiments are presented in Appendix A.2.

4.3 Ablation Studies

Regression models and aggregation approaches. Detailed results with various regression models and

UQ Method	XSUM		PubMedQA		CoQA		Mean Rank
	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	Acc.	AlignSc.	
MSP	-.356	-.153	-.024	.033	.648	.557	5.33
Focus	-.356	-.110	.045	-.063	.336	.261	6.50
SAR	-.029	.038	.075	.012	.474	.489	5.17
TAD (LinReg)	<u>.358</u>	<u>.223</u>	.429	.220	.639	.561	2.17
TAD (LinReg+AlignSc.)	.579	.345	-.018	.083	.657	.567	2.67
Gen. TAD (LinReg)	.006	-.032	.256	.208	.672	.541	3.33
Gen. TAD (LinReg+AlignSc.)	.210	.108	.179	.096	.675	.547	2.83

Table 3: The comparison of TAD trained on in-domain data with TAD trained on all out-of-domain datasets (designated with “Gen.”) (PRR \uparrow , Llama 8b v3 model). Warmer colors indicate better results. The best method is in bold, the second best is underlined.

aggregation approaches are presented in Table 4 and in Tables 12 and 13 in Appendix A. The optimal values of the hyper-parameters of TAD for all experimental setups are presented in Tables 14 to 16 in Appendix C.1 for Gemma 7b, LLaMA 8b v3, and StableLM 12b v2 models, respectively.

The results show that TAD based on regression using MLP and LinReg consistently outperform TAD based on CatBoost (Prokhorenkova et al., 2018). However, there is no big difference between MLP and LinReg. Therefore, for simplicity, we use LinReg as a regression method for TAD.

We investigate two strategies for aggregation of token-level TAD scores: the mean of the scores and the sum of the log scores inspired by perplexity. For the majority of the considered settings, the mean of the probabilities yields the best results. However, for QA with short answers, the sum of the log probabilities performs slightly better.

Comparison of features. Table 6 presents the experiments with various features for the regression model. For “TAD Embeds.”, we utilize the embeddings from the last hidden state from the decoder. For “TAD Probs.”, we use only generated probabilities for current and previous tokens, and $p(s_{i-1} = T)$. For “TAD Attn. Only”, we use attention, but without probabilities. TAD trained on attention weights with probabilities substantially outperforms all other options. We also note that TAD trained only on embeddings performs is much worse than other versions, emphasizing the importance of usage attention and probabilities.

Comparison to directly learning the unconditional probability. Table 5 compares TAD to directly learning the unconditional probability, where instead of using the target from equation 2, we simply try to approximate $p(s_i = T)$. These results demonstrate that the attention weights contain a lot of information about the unconditional probability

UQ Method	Aggregation	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
		ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
TAD (CatBoost)	$\frac{1}{K} \sum_{k=1}^K p_k$.496	.215	.201	.248	.064	-.011	.540	.181	.792	.382	.414	.283	.632	.578	.687	.634	.816	.800	5.89
TAD (CatBoost+AlignScore)	$\frac{1}{K} \sum_{k=1}^K p_k$.332	.146	.211	.269	.052	-.012	.556	.215	.665	.357	.382	.210	.603	.550	.550	.529	.818	.801	6.67
TAD (CatBoost)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.324	.284	.100	.075	-.078	.107	-.373	-.112	-.461	.011	.452	.163	.669	.609	.810	.736	.792	.776	7.33
TAD (CatBoost+AlignScore)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.249	.297	.087	.039	-.169	.093	-.373	-.190	-.472	-.001	.310	.083	.717	.626	.830	.774	.789	.775	8.28
TAD (LinReg)	$\frac{1}{K} \sum_{k=1}^K p_k$.502	.257	.329	.263	.172	.078	.576	.242	.787	.376	.563	.294	.510	.488	.619	.585	.811	.789	5.39
TAD (LinReg+AlignScore)	$\frac{1}{K} \sum_{k=1}^K p_k$.541	.380	.353	.349	.146	.092	.007	.064	.491	.472	.505	.368	.471	.441	.484	.462	.805	.782	5.17
TAD (LinReg)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.396	.319	.072	.090	-.029	.092	-.387	-.116	-.460	.012	.373	.224	.671	.608	.820	.751	.782	.760	7.22
TAD (LinReg+AlignScore)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.373	.351	.176	.121	-.099	.101	-.569	-.198	-.473	.000	.430	.187	.671	.600	.834	.777	.784	.766	7.22
TAD (MLP)	$\frac{1}{K} \sum_{k=1}^K p_k$.504	.249	.246	.210	.180	.080	.564	.217	.794	.369	.577	.298	.665	.605	.686	.641	.813	.794	4.50
TAD (MLP+AlignScore)	$\frac{1}{K} \sum_{k=1}^K p_k$.536	.349	.321	.322	.118	.092	-.059	.021	.624	.418	.419	.298	.614	.559	.608	.590	.804	.781	5.56
TAD (MLP)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.380	.301	.082	.042	-.020	.090	-.359	-.112	-.461	.010	.509	.183	.675	.613	.821	.754	.787	.764	7.28
TAD (MLP+AlignScore)	$\frac{1}{K} \sum_{k=1}^K \log p_k$.363	.340	.162	.105	-.100	.099	-.567	-.199	-.474	-.001	.220	.050	.713	.629	.836	.780	.789	.770	7.50

Table 4: Comparison of various considered regression models and aggregation strategies for TAD (PRR \uparrow , Gemma 7b model). Warmer colors indicate better results. The best method is in bold, the second best is underlined.

UQ Method	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Median Rel. Impr.
	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	Acc.	AlignSc.	Acc.	AlignSc.	Acc.	AlignSc.	Acc.	AlignSc.	
Learning $p(s_i = T)$.526	.345	.279	.314	.182	.079	.014	.015	.577	.471	.460	.389	.657	.591	.809	.774	.743	.760	-
TAD (LinReg+AlignSc.)	.541	.380	.353	.349	.146	.092	.007	.064	.491	.472	.505	.368	.471	.441	.484	.462	.805	.782	+3.1%

Table 5: The comparison of TAD with directly learning the unconditional probability $p(s_i = T)$ (PRR \uparrow , Gemma 7b model). The best method is in bold, the second best is underlined.

UQ Method	XSUM		PubMedQA		CoQA		Mean Rank
	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	Acc.	AlignSc.	
MSP	.329	.116	.455	.154	.699	.626	3.83
TAD Embeds. (LinReg. +AlignScore)	.191	.070	.025	.015	.606	.548	3.50
TAD Probs. (LinReg. +AlignScore)	.265	.234	.360	.142	.712	.613	2.83
TAD Attm. Only (LinReg+AlignScore)	.369	.252	.345	.112	.675	.608	2.67
TAD (LinReg+AlignScore)	.541	.380	.007	.064	.671	.600	2.17

Table 6: The comparison of various features for TAD (PRR \uparrow , Gemma 7b model). The best method is in bold, the second best is underlined.

UQ Method	Runtime per batch	Overhead
MSP	10.26 \pm 2.78	—
Mean Token Entropy	10.29 \pm 2.79	0.26%
Focus	10.55 \pm 2.84	2.80%
EigValLaplacian	44.90 \pm 9.55	340%
MC SE	44.72 \pm 9.53	340%
Semantic Entropy	44.87 \pm 9.54	340%
SAR	57.63 \pm 12.57	460%
TAD (CatBoost)	10.34 \pm 2.80	0.80%
TAD (LinReg)	10.27 \pm 2.78	0.10%
TAD (MLP)	10.27 \pm 2.78	0.11%

Table 7: The evaluation of the runtime of UQ methods measured on 900 instances from all datasets with predictions from Llama 8b v3. The best results are in bold.

itself. Nevertheless, TAD’s superior results show that taking into account the conditional dependency on previous generation steps and their uncertainty is also important.

4.4 Computational Efficiency

To demonstrate the computational efficiency of TAD, we compare its runtime to other UQ methods. We conducted experiments on 100 randomly sampled texts from each of our nine evaluation datasets using the LLaMA 8b v3 model on a single 80GB A100 GPU. The inference is implemented as a single-batch model call for all tokens in the output text. We use the LM-Polygraph (Fadeeva et al., 2023) implementation for other UQ methods.

Table 7 presents the average runtime per text

sample 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 (black-box, MC SE, Semantic Entropy, SAR) introduce huge computational overhead (340-460%) because they need to perform sampling from the LLM multiple times. On the contrary, TAD introduces minimal overhead (0.1-0.8%), which is much more practical.

5 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 on subsequent generation steps. This yields improved results in selective generation tasks, especially when the LLM output is long. Our experimental study shows that our proposed technique usually outperforms other state-of-the-art UQ methods (such as SAR) resulting in the best overall performance across three LLMs and nine datasets. TAD does not introduce much computational overhead due to the simplicity of the regression model (linear regression), which makes it a potentially practical choice for LLM-based applications.

In future work, we aim to apply the suggested method to quantifying the uncertainty of retrieval-augmented LLMs. TAD potentially could be used to take into account the credibility of the retrieved evidence.

618 Limitations

619 In the motivation of our approach, we assume a
620 strict Markov chain property between the genera-
621 tion steps. However, in reality, this property does
622 not hold as the current generation step usually de-
623 pends on multiple previous steps. This limitation
624 of our method could be addressed by estimating the
625 conditional dependency between multiple previous
626 steps, e.g., by using a Transformer layer instead
627 of the linear regressor. Nevertheless, our current
628 implementation that makes the Markov assumption
629 already yields strong results, and thus we leave
630 investigation of more complex modifications for
631 future work.

632 We also did not test our method on extra large
633 LLMs such as LLaMA 3 70b. We only used 7-12b
634 models due to limitations in our available computa-
635 tional resources.

636 Ethical Considerations

637 In our work, we considered open-source LLMs and
638 datasets not aimed at harmful content. However,
639 LLMs may generate potentially damaging texts for
640 various groups of people. Uncertainty quantifica-
641 tion techniques can help create more reliable use
642 of neural networks. Moreover, they can be applied
643 to detecting harmful generation, but this is not our
644 intention.

645 Moreover, despite that our proposed method
646 demonstrates significant performance improve-
647 ments, it can still mistakenly highlight correct and
648 not dangerous generated text with high uncertainty
649 in some cases. Thus, as with other uncertainty
650 quantification methods, it has limited application
651 for various tasks.

652 References

653 Asma Ben Abacha and Dina Demner-Fushman. 2019. [A](#)
654 [question-entailment approach to question answering](#).
655 *BMC Bioinform.*, 20(1):511:1–511:23.

656 Amos Azaria and Tom Mitchell. 2023. [The internal](#)
657 [state of an LLM knows when it’s lying](#). In *Find-*
658 *ings of the Association for Computational Linguistics:*
659 *EMNLP 2023*, pages 967–976, Singapore. Associa-
660 tion for Computational Linguistics.

661 Joris Baan, Nico Daheim, Evgenia Ilia, Dennis Ul-
662 mer, Haau-Sing Li, Raquel Fernández, Barbara
663 Plank, Rico Sennrich, Chrysoula Zerva, and Wilker
664 Aziz. 2023. Uncertainty in natural language gener-
665 ation: From theory to applications. *arXiv preprint*
666 *arXiv:2307.15703*.

667 Marco Bellagente, Jonathan Tow, Dakota Mahan, Duy
668 Phung, Maksym Zhuravinskyi, Reshinh Adithyan,
669 James Baicoianu, Ben Brooks, Nathan Cooper,
670 Ashish Datta, et al. 2024. Stable lm 2 1.6 b tech-
671 nical report. *arXiv preprint arXiv:2402.17834*.

672 Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen,
673 Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li,
674 and Yanghua Xiao. 2023. Hallucination detection:
675 Robustly discerning reliable answers in large lan-
676 guage models. In *Proceedings of the 32nd ACM*
677 *International Conference on Information and Knowl-*
678 *edge Management*, pages 245–255.

679 Julius Cheng and Andreas Vlachos. 2024. [Measur-](#)
680 [ing uncertainty in neural machine translation with](#)
681 [similarity-sensitive entropy](#). In *Proceedings of the*
682 *18th Conference of the European Chapter of the As-*
683 *sociation for Computational Linguistics (Volume 1:*
684 *Long Papers)*, pages 2115–2128, St. Julian’s, Malta.
685 Association for Computational Linguistics.

686 Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny,
687 Chenan Wang, Renjing Xu, Bhavya Kailkhura, and
688 Kaidi Xu. 2023. [Shifting attention to relevance: To-](#)
689 [wards the uncertainty estimation of large language](#)
690 [models](#). *Preprint*, arXiv:2307.01379.

691 Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem
692 Shelmanov, Sergey Petrakov, Haonan Li, Hamdy
693 Mubarak, Evgenii Tsybalov, Gleb Kuzmin, Alexan-
694 der Panchenko, Timothy Baldwin, et al. 2024. [Fact-](#)
695 [checking the output of large language models via](#)
696 [token-level uncertainty quantification](#). *arXiv preprint*
697 *arXiv:2403.04696*.

698 Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun,
699 Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin,
700 Daniil Vasilev, Elizaveta Goncharova, Alexander
701 Panchenko, Maxim Panov, Timothy Baldwin, and
702 Artem Shelmanov. 2023. [LM-polygraph: Uncer-](#)
703 [tainty estimation for language models](#). In *Proceed-*
704 *ings of the 2023 Conference on Empirical Methods*
705 *in Natural Language Processing: System Demon-*
706 *strations*, pages 446–461, Singapore. Association for
707 Computational Linguistics.

708 Marina Fomicheva, Shuo Sun, Lisa Yankovskaya,
709 Frédéric Blain, Francisco Guzmán, Mark Fishel,
710 Nikolaos Aletras, Vishrav Chaudhary, and Lucia Spe-
711 cia. 2020. [Unsupervised quality estimation for neural](#)
712 [machine translation](#). *Transactions of the Association*
713 *for Computational Linguistics*, 8:539–555.

714 Yarín Gal and Zoubin Ghahramani. 2016. [Dropout as](#)
715 [a Bayesian approximation: Representing model un-](#)
716 [certainty in deep learning](#). In *Proceedings of The*
717 *33rd International Conference on Machine Learn-*
718 *ing*, volume 48 of *Proceedings of Machine Learning*
719 *Research*, pages 1050–1059, New York, New York,
720 USA. PMLR.

721 Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppel,
722 Preslav Nakov, and Iryna Gurevych. 2023. A sur-
723 vey of language model confidence estimation and
724 calibration. *arXiv preprint arXiv:2311.08298*.

725	Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization . In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> , pages 70–79, Hong Kong, China. Association for Computational Linguistics.	783
726		784
727		785
728		786
729		
730		
731		
732	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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8362–8372.	
733		
734		
735		
736		
737		
738		
739	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A dataset for biomedical research question answering . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2567–2577, Hong Kong, China. Association for Computational Linguistics.	
740		
741		
742		
743		
744		
745		
746		
747		
748	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.	
749		
750		
751		
752		
753		
754		
755	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. <i>arXiv preprint arXiv:2207.05221</i> .	
756		
757		
758		
759		
760		
761	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
762		
763		
764		
765		
766		
767	Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym Korablyov, and Yoshua Bengio. 2022. DEUP: Direct epistemic uncertainty prediction. <i>Transactions on Machine Learning Research</i> .	
768		
769		
770		
771		
772	Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	
773		
774		
775		
776		
777	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	
778		
779		
780		
781		
782		
	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models . <i>CoRR</i> , abs/2305.19187.	783
		784
		785
		786
	Yu Lu, Jiali Zeng, Jiajun Zhang, Shuangzhi Wu, and Mu Li. 2022. Learning confidence for transformer-based neural machine translation . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2353–2364, Dublin, Ireland. Association for Computational Linguistics.	787
		788
		789
		790
		791
		792
		793
	Andrey Malinin and Mark J. F. Gales. 2021. Uncertainty estimation in autoregressive structured prediction . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	794
		795
		796
		797
		798
	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9004–9017, Singapore. Association for Computational Linguistics.	799
		800
		801
		802
		803
		804
		805
	Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, and et al. 2024. Gemma: Open models based on Gemini research and technology . <i>CoRR</i> , abs/2403.08295.	806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
	Sewon Min, Kalpesh Krishna, Xinxu 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 <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100.	824
		825
		826
		827
		828
		829
		830
	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 <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	831
		832
		833
		834
		835
		836
		837
	Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Fine-grained uncertainty quantification for	838
		839
		840

841	Ilms from semantic similarities. <i>arXiv preprint arXiv:2405.20003</i> .	In <i>Proceedings of the 3rd Workshop on Noisy User-generated Text</i> , pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.	896
842			897
843	Yookoon Park and David Blei. 2024. Density uncertainty layers for reliable uncertainty estimation. In <i>International Conference on Artificial Intelligence and Statistics</i> , pages 163–171. PMLR.	Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. The art of abstention: Selective prediction and error regularization for natural language processing . In <i>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)</i> , pages 1040–1051, Online. Association for Computational Linguistics.	899
844			900
845			901
846			902
847	Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. 2018. Catboost: unbiased boosting with categorical features. <i>Advances in neural information processing systems</i> , 31.		903
848			904
849			905
850			906
851			907
852	Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge . <i>Transactions of the Association for Computational Linguistics</i> , 7:249–266.	Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. AlignScore: Evaluating factual consistency with a unified alignment function . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.	908
853			909
854			910
855			911
856	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 <i>The Eleventh International Conference on Learning Representations</i> .		912
857			913
858			914
859			915
860			916
861			917
862	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. 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 <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 915–932, Singapore. Association for Computational Linguistics.	918
863			919
864			920
865			921
866			922
867			923
868			924
869	Junya Takayama and Yuki Arase. 2019. Relevant and informative response generation using pointwise mutual information . In <i>Proceedings of the First Workshop on NLP for Conversational AI</i> , pages 133–138. Association for Computational Linguistics.	Xuchao Zhang, Fanglan Chen, Chang-Tien Lu, and Naren Ramakrishnan. 2019. Mitigating uncertainty in document classification . In <i>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)</i> , pages 3126–3136, Minneapolis, Minnesota. Association for Computational Linguistics.	925
870			926
871			927
872			928
873			929
874	Liam van der Poel, Ryan Cotterell, and Clara Meister. 2022. Mutual information alleviates hallucinations in abstractive summarization . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5956–5965. Association for Computational Linguistics.		930
875			
876			
877			
878			
879			
880	Artem Vazhentsev, Gleb Kuzmin, Akim Tsvigun, Alexander Panchenko, Maxim Panov, Mikhail Burtsev, and Artem Shelmanov. 2023. Hybrid uncertainty quantification for selective text classification in ambiguous tasks . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11659–11681, Toronto, Canada. Association for Computational Linguistics.		
881			
882			
883			
884			
885			
886			
887			
888			
889	Yuxia Wang, Daniel Beck, Timothy Baldwin, and Karin Verspoor. 2022. Uncertainty estimation and reduction of pre-trained models for text regression . <i>Transactions of the Association for Computational Linguistics</i> , 10:680–696.		
890			
891			
892			
893			
894	Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions .		
895			

A Additional Experimental Results

A.1 Comparison with other UQ Methods

Here, we present the main results with Llama and StableLM.

UQ Method	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
MSP	<u>-386</u>	<u>-153</u>	<u>358</u>	<u>133</u>	<u>002</u>	<u>027</u>	<u>-024</u>	<u>033</u>	<u>417</u>	<u>493</u>	<u>324</u>	<u>174</u>	<u>648</u>	<u>557</u>	<u>671</u>	<u>590</u>	<u>252</u>	<u>706</u>	7.17
Perplexity	-388	-124	-088	231	130	196	222	-023	489	513	166	129	439	413	-456	-457	749	696	8.44
Mean Token Entropy	-385	-124	-114	230	132	189	233	-035	489	509	122	119	350	353	-498	-481	756	708	9.11
Focus	-356	-110	-024	<u>253</u>	112	<u>201</u>	045	-063	554	540	262	<u>274</u>	336	261	-469	-377	586	587	9.50
NumSemSets	011	062	154	185	070	099	005	037	-022	098	032	168	146	288	154	232	563	657	10.78
DegMat	048	085	191	215	076	100	013	027	069	174	112	145	306	440	317	405	633	697	8.83
Eccentricity	-009	036	034	<u>073</u>	042	054	-012	-008	048	062	086	046	484	476	386	443	643	652	11.33
EigValLaplacian	050	086	183	217	081	100	004	029	063	172	137	166	436	478	388	450	638	687	7.94
Lexical Similarity	011	038	302	182	105	093	099	025	272	143	012	012	482	473	372	414	652	647	9.22
MC NSE	-058	006	216	167	117	083	070	-006	304	217	013	012	441	407	038	071	656	657	10.72
MC SE	029	024	253	151	071	048	029	017	101	019	134	024	511	446	425	432	633	618	10.22
Semantic Entropy	029	026	256	157	066	050	031	015	102	022	121	023	521	483	444	459	686	675	9.06
SentenceSAR	-095	-005	167	125	053	033	-028	000	033	106	203	091	584	531	547	517	729	715	9.39
SAR	-029	038	288	208	115	112	075	012	328	237	012	085	474	489	149	181	718	721	7.78
TAD (LinReg)	<u>358</u>	<u>223</u>	<u>336</u>	<u>219</u>	<u>210</u>	<u>111</u>	<u>429</u>	<u>220</u>	500	501	189	130	639	561	868	758	707	671	3.32
TAD (LinReg+AlignScore)	<u>579</u>	<u>345</u>	<u>404</u>	<u>369</u>	207	150	-018	083	613	544	251	235	657	567	914	824	715	691	2.78

Table 8: PRR \uparrow of UQ methods for the Llama 8b v3 model. Warmer colors indicate better results. The best method is in bold, the second best is underlined.

UQ Method	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
MSP	-144	-060	498	341	-027	062	-429	-168	478	596	<u>250</u>	<u>224</u>	680	597	717	685	738	715	6.83
Perplexity	-257	-034	434	351	092	044	409	099	492	592	219	179	385	373	-340	-385	732	700	8.33
Mean Token Entropy	-250	-028	409	340	108	034	410	085	503	593	139	181	312	329	-403	-423	747	713	9.00
Focus	-173	019	300	228	040	011	214	080	552	639	217	215	147	105	-165	-197	643	649	10.50
NumSemSets	001	054	179	187	005	074	081	051	-007	055	060	167	221	303	110	200	576	636	12.83
DegMat	-000	057	309	326	017	120	052	039	136	242	214	194	342	489	452	561	653	698	9.33
Eccentricity	-034	004	238	250	023	049	-025	007	146	179	165	047	527	557	496	568	643	660	11.61
EigValLaplacian	-008	063	292	311	012	115	049	038	116	226	227	215	500	557	513	581	661	697	8.78
Lexical Similarity	-111	<u>079</u>	381	285	-119	098	094	026	296	271	141	090	508	524	489	545	656	670	8.56
MC NSE	068	048	371	263	073	088	161	059	370	372	123	126	437	421	273	310	623	615	10.22
MC SE	066	-006	393	291	059	068	034	026	209	234	164	051	565	527	515	537	623	616	10.67
Semantic Entropy	067	-003	412	317	066	071	033	024	215	247	152	047	578	565	545	578	674	670	9.06
SentenceSAR	005	001	392	330	010	044	-052	001	255	307	280	157	642	603	630	644	713	713	8.22
SAR	079	079	412	341	080	119	177	059	405	401	209	196	494	531	398	460	702	714	6.33
TAD (LinReg)	<u>375</u>	<u>224</u>	<u>459</u>	<u>282</u>	<u>163</u>	<u>137</u>	<u>493</u>	<u>284</u>	511	610	368	222	707	624	850	786	688	671	3.50
TAD (LinReg+AlignScore)	<u>459</u>	068	<u>519</u>	<u>419</u>	145	122	249	212	696	674	462	367	658	614	863	803	696	691	2.22

Table 9: PRR \uparrow of UQ methods for the StableLM 12b v2 model. Warmer colors indicate better results. The best method is in bold, the second best is underlined.

A.2 Generalization Experiments

Tables 3, 10 and 11 present the comparison of the TAD trained on the in-domain training dataset with the TAD trained on all out-of-domain datasets for Gemma 7b, Llama 8b v3, and StableLM 12b v2 models respectively. In this experiment, we examine how our approach can be generalized on the unseen datasets. For each dataset, we create a general training dataset by using 300 samples from the training datasets from each of the eight other datasets used in the experiments. Thus, we evaluate TAD that is not trained on the target dataset. We conduct experiments on one dataset from each task: XSUM, PubMedQA, and CoQA. We compare the results with three strongest baseline methods: MSP, Focus, and SAR. Overall, we can see that the TAD method can be generalized on the unseen datasets and outperform all other baselines in most settings.

UQ Method	XSUM		PubMedQA		CoQA		Mean Rank
	ROUGE-L	AlignSc.	ROUGE-L	AlignSc.	Acc.	AlignSc.	
MSP	-0.329	-0.116	-0.455	-0.154	0.699	0.626	5.00
Focus	-0.324	-0.161	-0.357	-0.146	0.322	0.250	6.50
SAR	0.042	-0.006	0.111	0.014	0.477	0.503	4.50
TAD (LinReg)	0.502	0.257	0.576	0.242	0.671	0.608	2.17
TAD (LinReg+AlignSc.)	0.541	0.380	0.07	0.064	0.671	0.600	2.67
Gen. TAD (LinReg)	-0.061	-0.068	0.288	0.101	0.703	0.594	3.17
Gen. TAD (LinReg+AlignSc.)	0.132	0.096	-0.124	-0.074	0.696	0.589	4.00

Table 10: The comparison of TAD trained on in-domain data with TAD trained on all out-of-domain datasets (designated with “Gen.”) (PRR \uparrow , Gemma 7b model). Warmer colors indicate better results. The best method is in bold, the second best is underlined.

UQ Method	XSUM		PubMedQA		CoQA		Mean Rank
	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	
MSP	<u>-144</u>	<u>-060</u>	<u>-429</u>	<u>-168</u>	<u>.680</u>	<u>.597</u>	5.83
Focus	<u>-173</u>	<u>.019</u>	<u>.214</u>	<u>.080</u>	<u>.147</u>	<u>.105</u>	5.83
SAR	<u>.079</u>	.079	<u>.177</u>	<u>.059</u>	<u>.494</u>	<u>.531</u>	4.67
TAD (LinReg)	<u>.375</u>	<u>.024</u>	.493	.284	<u>.707</u>	<u>.624</u>	1.67
TAD (LinReg, +AlignScore)	.459	<u>.068</u>	<u>.249</u>	<u>.219</u>	<u>.698</u>	<u>.614</u>	<u>2.50</u>
Gen. TAD (LinReg)	<u>-032</u>	<u>-015</u>	<u>.433</u>	<u>.217</u>	<u>.701</u>	<u>.584</u>	4.00
Gen. TAD (LinReg, +AlignScore)	<u>.023</u>	<u>-008</u>	<u>.288</u>	<u>.143</u>	.709	<u>.592</u>	3.50

Table 11: The comparison of TAD trained on in-domain data with TAD trained on all out-of-domain datasets (designated with ‘‘Gen.’’) (PRR \uparrow , StableLM 12b v2 model). Warmer colors indicate better results. The best method is in bold, the second best is underlined.

A.3 Ablation Studies

Here, we present ablation studies for regression models and aggregation techniques with additional LLMs.

UQ Method	Aggregation	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
		ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
TAD (CatBoost)	$\frac{1}{k} \sum_{k=1}^k p_k$.349	.183	<u>-064</u>	.211	.180	.101	<u>.366</u>	.150	.448	.476	.208	.146	.605	.536	.741	.665	<u>.743</u>	.710	6.78
TAD (CatBoost, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$	<u>.250</u>	<u>.097</u>	<u>-013</u>	<u>.258</u>	<u>.137</u>	<u>.092</u>	<u>.255</u>	<u>.002</u>	.448	.492	.234	.179	.576	.509	<u>.805</u>	<u>.558</u>	.746	.714	6.94
TAD (CatBoost)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.357	.244	.279	.026	<u>-068</u>	<u>-036</u>	<u>-429</u>	<u>-056</u>	.293	.411	.323	.191	.647	.557	.813	.708	.715	.680	7.89
TAD (CatBoost, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.272</u>	<u>.227</u>	<u>.269</u>	<u>.015</u>	<u>-115</u>	<u>-050</u>	<u>-461</u>	<u>-070</u>	<u>.099</u>	<u>.239</u>	<u>.305</u>	<u>.189</u>	.672	.566	.875	.795	.712	.683	8.61
TAD (LinReg)	$\frac{1}{k} \sum_{k=1}^k p_k$	<u>.358</u>	<u>.223</u>	<u>.336</u>	<u>.219</u>	<u>.210</u>	<u>.111</u>	.429	.220	<u>.500</u>	<u>.501</u>	<u>.189</u>	<u>.130</u>	<u>.535</u>	<u>.507</u>	<u>.742</u>	<u>.671</u>	<u>.739</u>	<u>.702</u>	6.11
TAD (LinReg, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$.579	.345	.404	.369	<u>.207</u>	<u>.150</u>	<u>-018</u>	<u>.083</u>	.613	.544	<u>.251</u>	<u>.235</u>	<u>.509</u>	<u>.473</u>	<u>.637</u>	<u>.591</u>	<u>.738</u>	<u>.708</u>	4.89
TAD (LinReg)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.438	.291	.307	.082	.005	-.021	<u>-402</u>	<u>-049</u>	.310	.421	.396	<u>.261</u>	.639	.561	.868	.758	<u>.707</u>	<u>.671</u>	6.61
TAD (LinReg, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.466</u>	<u>.334</u>	<u>.367</u>	<u>.098</u>	<u>-040</u>	<u>-041</u>	<u>-447</u>	<u>-065</u>	<u>.175</u>	<u>.256</u>	<u>.273</u>	<u>.195</u>	<u>.657</u>	<u>.567</u>	.914	.824	<u>.715</u>	<u>.691</u>	6.11
TAD (MLP)	$\frac{1}{k} \sum_{k=1}^k p_k$.496	.256	.317	.221	.215	.119	<u>.408</u>	<u>.166</u>	<u>.509</u>	<u>.488</u>	<u>.189</u>	<u>.132</u>	<u>.587</u>	<u>.525</u>	<u>.751</u>	<u>.664</u>	<u>.738</u>	<u>.701</u>	5.72
TAD (MLP, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$.572	<u>.326</u>	<u>.303</u>	.346	<u>.206</u>	<u>.145</u>	<u>.294</u>	.251	<u>.563</u>	<u>.487</u>	<u>.255</u>	<u>.276</u>	<u>.551</u>	<u>.494</u>	<u>.675</u>	<u>.635</u>	<u>.739</u>	.712	4.89
TAD (MLP)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.448	.303	.310	.069	.008	-.021	<u>-419</u>	<u>-056</u>	<u>.301</u>	<u>.407</u>	<u>.355</u>	<u>.238</u>	<u>.646</u>	<u>.566</u>	<u>.879</u>	<u>.757</u>	<u>.718</u>	<u>.682</u>	6.28
TAD (MLP, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.435</u>	<u>.326</u>	<u>.352</u>	<u>.088</u>	<u>-052</u>	<u>-046</u>	<u>-453</u>	<u>-063</u>	<u>.153</u>	<u>.220</u>	<u>.191</u>	<u>.146</u>	<u>.662</u>	.575	<u>.912</u>	<u>.822</u>	<u>.717</u>	<u>.693</u>	7.17

Table 12: Comparison of various considered regression models and different aggregation strategies for TAD by PRR \uparrow for the Llama 8b v3 model for various tasks. Warmer colors indicate better results. The best method is in bold, the second best is underlined.

UQ Method	Aggregation	XSUM		SamSum		CNN		PubMedQA		MedQUAD		TruthfulQA		CoQA		SciQ		TriviaQA		Mean Rank
		ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	ROUGE-L	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	Acc.	AlignScore	
TAD (CatBoost)	$\frac{1}{k} \sum_{k=1}^k p_k$.374	.019	<u>.409</u>	.296	<u>.117</u>	<u>.071</u>	.495	.278	.500	.586	<u>.394</u>	<u>.242</u>	<u>.637</u>	.574	.710	.678	<u>.725</u>	<u>.701</u>	6.67
TAD (CatBoost, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$	<u>.262</u>	<u>.036</u>	.440	.311	<u>.076</u>	<u>.060</u>	<u>.295</u>	<u>.107</u>	.525	.589	<u>.418</u>	<u>.297</u>	<u>.585</u>	<u>.529</u>	<u>.676</u>	<u>.650</u>	.726	.703	6.94
TAD (CatBoost)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.320	<u>-029</u>	.442	.296	<u>-030</u>	<u>.151</u>	<u>-.565</u>	<u>-188</u>	<u>.452</u>	<u>.586</u>	.539	<u>.236</u>	<u>.703</u>	<u>.619</u>	<u>.826</u>	<u>.763</u>	<u>.710</u>	<u>.675</u>	7.39
TAD (CatBoost, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.248</u>	<u>-021</u>	<u>.380</u>	<u>.255</u>	<u>-105</u>	<u>.110</u>	<u>-.582</u>	<u>-192</u>	<u>.316</u>	<u>.410</u>	<u>.451</u>	<u>.199</u>	.715	.620	.862	.802	<u>.702</u>	<u>.681</u>	8.94
TAD (LinReg)	$\frac{1}{k} \sum_{k=1}^k p_k$.375	.024	<u>.459</u>	<u>.282</u>	<u>.163</u>	<u>.137</u>	.493	.284	.511	.610	<u>.368</u>	<u>.222</u>	<u>.594</u>	<u>.555</u>	<u>.734</u>	<u>.710</u>	<u>.712</u>	<u>.686</u>	5.83
TAD (LinReg, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$.459	.068	.419	.419	<u>.145</u>	<u>.127</u>	<u>.249</u>	<u>.219</u>	<u>.096</u>	.474	<u>.462</u>	<u>.267</u>	<u>.536</u>	<u>.488</u>	<u>.684</u>	<u>.661</u>	<u>.710</u>	<u>.693</u>	5.86
TAD (LinReg)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.368	.011	.450	.279	.013	<u>.154</u>	<u>-.556</u>	<u>-185</u>	<u>.463</u>	<u>.599</u>	<u>.300</u>	<u>.228</u>	<u>.207</u>	.624	<u>.850</u>	<u>.786</u>	<u>.688</u>	<u>.653</u>	6.50
TAD (LinReg, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.358</u>	<u>.036</u>	<u>.442</u>	<u>.324</u>	<u>-023</u>	<u>.135</u>	<u>-.567</u>	<u>-186</u>	<u>.429</u>	<u>.436</u>	<u>.453</u>	<u>.243</u>	<u>.698</u>	<u>.614</u>	.863	.803	<u>.696</u>	<u>.674</u>	6.89
TAD (MLP)	$\frac{1}{k} \sum_{k=1}^k p_k$.401	.018	.473	.301	.166	<u>.149</u>	<u>.488</u>	<u>.283</u>	.516	.606	<u>.397</u>	<u>.237</u>	<u>.605</u>	<u>.554</u>	<u>.728</u>	<u>.708</u>	<u>.711</u>	<u>.684</u>	5.61
TAD (MLP, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k p_k$.460	.065	<u>.500</u>	<u>.383</u>	<u>.135</u>	<u>.114</u>	<u>.249</u>	<u>.236</u>	.722	.662	.525	.414	<u>.528</u>	<u>.488</u>	<u>.733</u>	<u>.710</u>	<u>.705</u>	<u>.692</u>	4.67
TAD (MLP)	$\frac{1}{k} \sum_{k=1}^k \log p_k$.363	.009	.458	.292	.024	.169	<u>-.557</u>	<u>-186</u>	<u>.475</u>	<u>.600</u>	<u>.307</u>	<u>.259</u>	<u>.701</u>	<u>.624</u>	<u>.843</u>	<u>.779</u>	<u>.695</u>	<u>.663</u>	6.39
TAD (MLP, +AlignScore)	$\frac{1}{k} \sum_{k=1}^k \log p_k$	<u>.343</u>	<u>.032</u>	<u>.441</u>	<u>.325</u>	<u>-014</u>	<u>.136</u>	<u>-.571</u>	<u>-186</u>	<u>.404</u>	<u>.396</u>	<u>.450</u>	<u>.227</u>	<u>.706</u>	<u>.621</u>	<u>.848</u>	<u>.790</u>	<u>.703</u>	<u>.684</u>	7.11

Table 13: Comparison of various considered regression models and different aggregation strategies for TAD by PRR \uparrow for StableLM 12b v2 model for various tasks. Warmer colors indicate better results. The best method is in bold, the second best is underlined.

B Computational Resources

All experiments were conducted on a single NVIDIA A100 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 14 to 16 for Gemma 7b, Llama 8b v3, and StableLM 12b v2 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 CatBoost is the following:

Num. of trees: [100, 200];

959
960
961
962
963
964
965
966
967

Learning rate: [1e-1, 1e-2];

Tree depth: [3, 5].

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].

UQ Method	Aggregation	XSUM	SamSum	CNN	PubMedQA	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA
TAD (CatBoost)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 3	100, 0.01, 3	100, 0.01, 3	100, 0.01, 5	100, 0.1, 5	100, 0.01, 3	100, 0.01, 3	100, 0.01, 3	100, 0.01, 5
TAD (CatBoost, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 5	200, 0.01, 3	100, 0.01, 3	100, 0.1, 5	100, 0.01, 3	100, 0.01, 5	100, 0.01, 5	200, 0.1, 5	100, 0.01, 5
TAD (CatBoost)	$\sum_{k=1}^K \log P_k$	200, 0.1, 3	200, 0.1, 5	200, 0.1, 5	200, 0.01, 3	200, 0.1, 5	200, 0.1, 5	100, 0.1, 5	200, 0.1, 3	200, 0.1, 5
TAD (CatBoost, +AlignScore)	$\sum_{k=1}^K \log P_k$	200, 0.1, 5	200, 0.1, 5	100, 0.01, 3	100, 0.1, 5	100, 0.01, 5	200, 0.1, 3	100, 0.01, 5	100, 0.01, 3	100, 0.01, 5
TAD (LinReg)	$\frac{1}{R} \sum_{k=1}^K P_k$	1	10.0	0.01	1	10.0	0.0001	10.0	10.0	10.0
TAD (LinReg, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	0.01	0.001	0.001	0.1	0.001	0.01	10.0	1	10.0
TAD (LinReg)	$\sum_{k=1}^K \log P_k$	10.0	0.0001	0.0001	10.0	0.0001	0.01	10.0	1	1
TAD (LinReg, +AlignScore)	$\sum_{k=1}^K \log P_k$	0.01	0.001	0.0001	0.0001	0.001	0.001	10.0	1	1
TAD (MLP)	$\frac{1}{R} \sum_{k=1}^K P_k$	2, 30, 3e-05, 128	2, 10, 1e-05, 128	2, 30, 5e-05, 128	4, 10, 3e-05, 64	2, 10, 1e-05, 128	4, 30, 5e-05, 128	2, 10, 1e-05, 128	2, 10, 1e-05, 128	2, 10, 1e-05, 128
TAD (MLP, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	2, 30, 3e-05, 64	2, 30, 3e-05, 128	2, 30, 5e-05, 128	4, 10, 5e-05, 128	2, 10, 1e-05, 128	2, 10, 3e-05, 64	2, 10, 1e-05, 128	4, 10, 5e-05, 64	4, 10, 5e-05, 64
TAD (MLP)	$\sum_{k=1}^K \log P_k$	2, 20, 5e-05, 64	2, 10, 1e-05, 128	4, 30, 5e-05, 128	4, 10, 1e-05, 128	4, 30, 3e-05, 128	4, 20, 3e-05, 64	4, 10, 1e-05, 64	4, 30, 3e-05, 128	4, 30, 1e-05, 64
TAD (MLP, +AlignScore)	$\sum_{k=1}^K \log P_k$	4, 20, 5e-05, 128	4, 30, 5e-05, 128	4, 20, 5e-05, 128	4, 30, 5e-05, 64	4, 30, 5e-05, 64	4, 30, 5e-05, 128	2, 20, 1e-05, 128	4, 20, 3e-05, 128	4, 30, 1e-05, 64

Table 14: Optimal hyperparameters for the TAD methods for the Gemma 7b model.

UQ Method	Aggregation	XSUM	SamSum	CNN	PubMedQA	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA
TAD (CatBoost)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 5	100, 0.01, 3	200, 0.1, 5	100, 0.01, 3	100, 0.01, 3	200, 0.1, 5	100, 0.01, 5	100, 0.01, 3	100, 0.01, 3
TAD (CatBoost, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 5	200, 0.01, 3	100, 0.01, 3	100, 0.01, 5	200, 0.01, 5	200, 0.1, 5	100, 0.01, 5	200, 0.1, 5	100, 0.01, 5
TAD (CatBoost)	$\sum_{k=1}^K \log P_k$	200, 0.1, 5	100, 0.01, 3	200, 0.1, 5	200, 0.1, 5	200, 0.1, 5	200, 0.1, 5	200, 0.1, 3	200, 0.1, 3	200, 0.1, 3
TAD (CatBoost, +AlignScore)	$\sum_{k=1}^K \log P_k$	200, 0.1, 5	100, 0.01, 3	100, 0.01, 5	100, 0.01, 5	100, 0.01, 5	200, 0.1, 5	100, 0.1, 3	200, 0.1, 5	100, 0.01, 5
TAD (LinReg)	$\frac{1}{R} \sum_{k=1}^K P_k$	0.0001	10.0	0.01	0.1	0.0001	10.0	10.0	10.0	10.0
TAD (LinReg, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	0.001	0.0001	0.01	0.01	0.0001	0.1	10.0	1	10.0
TAD (LinReg)	$\sum_{k=1}^K \log P_k$	0.01	1	0.001	0.0001	0.0001	0.0001	10.0	1	10.0
TAD (LinReg, +AlignScore)	$\sum_{k=1}^K \log P_k$	0.0001	0.0001	0.0001	0.1	0.0001	0.1	10.0	1	10.0
TAD (MLP)	$\frac{1}{R} \sum_{k=1}^K P_k$	2, 10, 1e-05, 64	4, 30, 5e-05, 128	2, 30, 5e-05, 128	4, 10, 5e-05, 64	2, 20, 5e-05, 128	2, 30, 3e-05, 128	2, 10, 1e-05, 128	2, 30, 1e-05, 128	4, 30, 1e-05, 128
TAD (MLP, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	4, 10, 3e-05, 128	4, 20, 1e-05, 128	2, 20, 5e-05, 128	4, 10, 5e-05, 64	4, 30, 1e-05, 128	4, 30, 5e-05, 128	2, 10, 1e-05, 128	2, 20, 5e-05, 64	4, 10, 5e-05, 64
TAD (MLP)	$\sum_{k=1}^K \log P_k$	4, 10, 1e-05, 128	4, 10, 1e-05, 128	4, 30, 5e-05, 64	4, 20, 5e-05, 64	2, 30, 5e-05, 128	4, 30, 3e-05, 64	4, 10, 1e-05, 64	2, 10, 3e-05, 128	4, 20, 5e-05, 128
TAD (MLP, +AlignScore)	$\sum_{k=1}^K \log P_k$	2, 30, 1e-05, 128	4, 30, 3e-05, 64	2, 30, 5e-05, 64	2, 20, 5e-05, 64	4, 30, 1e-05, 128	4, 30, 3e-05, 128	2, 10, 3e-05, 128	2, 30, 3e-05, 128	4, 10, 5e-05, 128

Table 15: Optimal hyperparameters for the TAD methods for the Llama 8b v3 model.

UQ Method	Aggregation	XSUM	SamSum	CNN	PubMedQA	MedQUAD	TruthfulQA	CoQA	SciQ	TriviaQA
TAD (CatBoost)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 5	100, 0.01, 3	100, 0.01, 3	200, 0.1, 5	200, 0.1, 5	200, 0.1, 3	100, 0.01, 3	100, 0.01, 3	200, 0.1, 3
TAD (CatBoost, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	200, 0.1, 3	100, 0.01, 3	100, 0.01, 5	200, 0.01, 3	100, 0.1, 3	200, 0.1, 5	100, 0.01, 5	200, 0.1, 5	100, 0.01, 5
TAD (CatBoost)	$\sum_{k=1}^K \log P_k$	200, 0.1, 5	100, 0.1, 3	200, 0.1, 5	200, 0.1, 3	200, 0.1, 5	200, 0.1, 5	100, 0.1, 5	100, 0.01, 5	200, 0.1, 3
TAD (CatBoost, +AlignScore)	$\sum_{k=1}^K \log P_k$	200, 0.1, 3	100, 0.01, 3	100, 0.01, 5	200, 0.1, 5	100, 0.01, 5	200, 0.1, 5	100, 0.1, 3	100, 0.01, 5	200, 0.1, 3
TAD (LinReg)	$\frac{1}{R} \sum_{k=1}^K P_k$	0.01	10.0	1	10.0	0.0001	0.0001	10.0	10.0	10.0
TAD (LinReg, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	0.01	0.1	0.1	0.0001	0.001	0.1	10.0	10.0	10.0
TAD (LinReg)	$\sum_{k=1}^K \log P_k$	1	10.0	1	1	0.01	0.001	10.0	10.0	10.0
TAD (LinReg, +AlignScore)	$\sum_{k=1}^K \log P_k$	0.1	0.001	0.01	1	0.1	0.001	10.0	10.0	10.0
TAD (MLP)	$\frac{1}{R} \sum_{k=1}^K P_k$	4, 10, 1e-05, 128	4, 10, 5e-05, 64	2, 30, 1e-05, 128	2, 30, 1e-05, 128	4, 30, 3e-05, 64	4, 30, 1e-05, 128	2, 10, 1e-05, 128	4, 10, 3e-05, 64	4, 30, 1e-05, 128
TAD (MLP, +AlignScore)	$\frac{1}{R} \sum_{k=1}^K P_k$	2, 30, 3e-05, 128	4, 10, 1e-05, 128	4, 30, 5e-05, 64	4, 10, 1e-05, 128	4, 10, 5e-05, 128	4, 30, 1e-05, 64	2, 10, 3e-05, 128	4, 10, 3e-05, 64	4, 10, 3e-05, 64
TAD (MLP)	$\sum_{k=1}^K \log P_k$	2, 20, 1e-05, 128	4, 10, 5e-05, 64	4, 30, 5e-05, 64	4, 20, 5e-05, 64	4, 20, 5e-05, 64	4, 20, 5e-05, 64	4, 10, 3e-05, 64	2, 10, 1e-05, 64	4, 10, 3e-05, 64
TAD (MLP, +AlignScore)	$\sum_{k=1}^K \log P_k$	2, 30, 1e-05, 64	2, 30, 3e-05, 64	4, 30, 5e-05, 64	4, 20, 1e-05, 64	4, 30, 3e-05, 128	4, 30, 5e-05, 128	2, 10, 1e-05, 64	4, 10, 1e-05, 128	4, 10, 3e-05, 64

Table 16: Optimal hyperparameters for the TAD methods for the StableLM 12b v2 model.

C.2 LLM Generation Hyperparameters

968

Dataset	Task	Max Input Length	Generation Length	Temperature	Top-p	Do Sample	Beams	Repetition Penalty
XSum	TS		56					
SamSum			128					
CNN			128					
PubMedQA	QA Long answer	-	128	1.0	1.0	False	1	1
MedQUAD			128					
TruthfulQA			128					
CoQA	QA Short answer		20					
SciQ			20					
TriviQA			20					

Table 17: Text generation hyperparameters for all LLMs used in the experiments.

D Dataset Statistics

969

Task	Dataset	N-shot	Train texts for TAD	Evaluation texts
Text Summarization	CNN/DailyMail	0	2,000	2,000
	XSum	0	2,000	2,000
	SamSum	0	2,000	819
QA Long answer	PubMedQA	0	2,000	2,000
	MedQUAD	5	1,000	2,000
	TruthfulQA	5	408	409
QA Short answer	SciQ	0	2,000	1,000
	CoQA	all preceding questions	2,000	2,000
	TriviQA	5	2,000	2,000

Table 18: The statistics of the datasets used for evaluation.