

# UALIGN: Leveraging Uncertainty Estimations for Factuality Alignment on Large Language Models

Anonymous ACL submission

## Abstract

Despite demonstrating impressive capabilities, Large Language Models (LLMs) still often struggle to accurately express the factual knowledge they possess, especially in cases where the LLMs’ knowledge boundaries are ambiguous. To improve LLMs’ factual expressions, we propose the UALIGN framework, which leverages Uncertainty estimations to represent knowledge boundaries, and then explicitly incorporates these representations as input features into prompts for LLMs to **Align** with factual knowledge. First, we prepare the dataset on knowledge question-answering (QA) samples by calculating two uncertainty estimations, including confidence score and semantic entropy, to represent the knowledge boundaries for LLMs. Subsequently, using the prepared dataset, we train a reward model that incorporates uncertainty estimations and then employ the Proximal Policy Optimization (PPO) algorithm for factuality alignment on LLMs. Experimental results indicate that, by integrating uncertainty representations in LLM alignment, the proposed UALIGN can significantly enhance the LLMs’ capacities to confidently answer known questions and refuse unknown questions on both in-domain and out-of-domain tasks, showing reliability improvements and good generalizability over various prompt- and training-based baselines.

## 1 Introduction

Despite the remarkable proficiency of large language models (LLMs) across a diverse range of tasks (Touvron et al., 2023; OpenAI, 2023; Chiang et al., 2023), they still frequently face challenges in accurately expressing factual knowledge that they learned from the pre-training stage but are uncertain about. In such cases, the knowledge boundaries are somewhat ambiguous by LLMs, remaining a gap between “known” and “expression” (Lin et al., 2024; Zhang et al., 2024b; Li et al., 2024),

which may lead to the hallucination problem and undermine the reliability and applicability to users.

LLMs typically generate responses (“expression”) based on knowledge distributions learned during pre-training (“known”). However, much of the knowledge acquired during this phase exhibits vague boundaries, comprising numerous learned but uncertain knowledge pieces (*weakly known*, *light green area of spectrum* in Fig. 1 (a)) (Gekhman et al., 2024). Hence, LLMs may not confidently convey accurate information in downstream tasks even though they hold relevant knowledge but don’t make sure (Zhang et al., 2024b). Additionally, LLMs may exhibit overconfidence in the knowledge they are unfamiliar with (*unknown*, *the gray area of spectrum* in Fig. 1 (a)), leading to fabricated or hallucinatory content (Zhang et al., 2024a; Liu et al., 2024). This issue primarily arises from that LLMs don’t properly reconcile the knowledge boundaries with factual accuracy during alignment (Tian et al., 2024). Unlike previous works that focused on reinforcement learning (RL) through knowledge feedback or factuality alignment (Liang et al., 2024; Xu et al., 2024a; Tian et al., 2024; Lin et al., 2024; Zhang et al., 2024b; Yang et al., 2024), our objective is to elicit LLMs’ weakly known facts and extend beyond merely discerning unknown facts by explicitly utilizing knowledge boundaries in alignment. We aim to leverage the knowledge boundary information of LLMs to instruct LLMs to confidently express their known yet uncertain information and firmly refuse questions beyond their knowledge as in Fig. 1 (b). Based on improvements of “known”, LLMs’ expressions are more truthful and reliable, thereby minimizing the discrepancy between “known” and “expression” (Lin et al., 2024; Zhang et al., 2024b; Li et al., 2024).

Inspired by the aforementioned analysis, we propose the UALIGN framework, which strategically models Uncertainty regarding knowledge boundary representations, subsequently **Aligning** these esti-

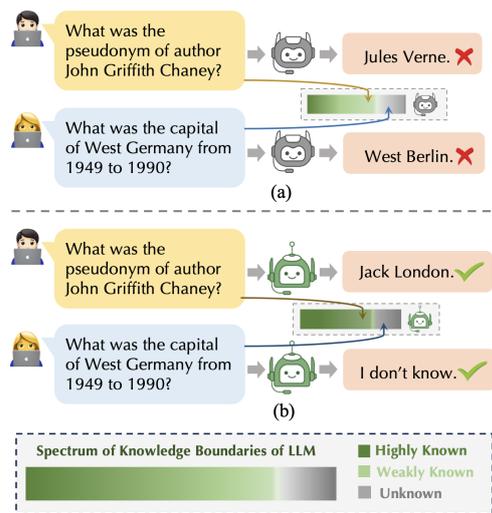


Figure 1: Examples of LLMs with (a) ambiguous and (b) explicit knowledge boundaries to answer questions.

mations with factuality. Therefore, the UALIGN framework focuses on two pivotal issues: how to capture the knowledge boundary representations and how to align with factuality.

First, we prepare the dataset that incorporates knowledge boundary information for alignment in the UALIGN framework. Knowledge boundaries always indicate the known level of factual knowledge, generally implemented using uncertainty estimation methods on LLMs (Ren et al., 2023). To precisely capture the intrinsic perception of knowledge boundary representations given the knowledge QA datasets, we adopt two uncertainty estimations of accuracy-based confidence score (Xiong et al., 2024) and semantic entropy (Kuhn et al., 2023) respectively. We sample multiple responses to a question using varied prompting and temperature sampling to approximate actual knowledge boundaries by calculating the confidence and entropy of each question. The two measures (Kuhn et al., 2023; Xiong et al., 2024), as complementary, can reflect the convince and dispersion of generated responses to a question based on LLMs’ internal knowledge. Questions with at least one correct sampled answer are regarded as “known”, and those with all incorrect sampled responses are considered “unknown”. We revise ground-truth answers to unknown questions to refusal responses to delineate known and unknown facts (Zhang et al., 2024a).

Second, following Ouyang et al. (2022), we explicitly leverage the uncertainty estimations to align with factuality on the prepared dataset using both supervised fine-tuning (SFT) and reinforcement learning (RL). We employ SFT to train two uncer-

tainty estimation models to predict confidence and entropy, and then train a reward model to evaluate the correctness of the generated answer conditioned on the input comprising the question, the generated response, and two uncertainty estimations regarding the knowledge boundary. With the reward model, we further adopt the Proximal Policy Optimization (PPO) (Schulman et al., 2017) algorithm for LLM alignment by feeding both questions and two measures as prompts to elicit the policy LLM’s factual expressions to improve the reliability.

Experiments are conducted to evaluate in-domain and out-of-domain performance on a range of knowledge QA datasets. The results demonstrate our proposed UALIGN method significantly enhances the reliability and generalization for LLMs over several baseline methods to accurately express known factual knowledge and refuse unknown questions, suggesting that leveraging the two employed uncertainty estimations in alignment can notably improve LLMs’ factuality.

In summary, our contributions are as follows.

1) To the best of our knowledge, UALIGN is the first to explicitly leverage the uncertainty estimations representing knowledge boundaries for LLM alignment, heralding a promising direction for future research of LLM training<sup>1</sup>.

2) We demonstrate that jointly incorporating confidence and semantic entropy into prompts can provide precise knowledge boundary information to elicit LLMs’ factual expressions.

3) We conduct main experiments by comparing our UALIGN with various baselines as well as ablation studies, validating the reliability improvements and robust generalization of the UALIGN method.

## 2 Methodology

The proposed UALIGN framework is introduced in this section with two parts: The Sec. 2.1 involves the UALIGN dataset preparation process, including strategies to collect multiple responses, as well as uncertainty measures to capture intrinsic representations of knowledge boundary on knowledge-based QA pairs as illustrated in Fig. 2. The Sec. 2.2 utilizes the obtained UALIGN dataset to train the uncertainty estimation models, and further explicitly incorporate the estimations as input features to elicit LLMs to generate factual responses using SFT- and PPO-based alignment methods as shown in Fig. 3 and Algorithm 1.

<sup>1</sup>The codes will be released on GitHub.

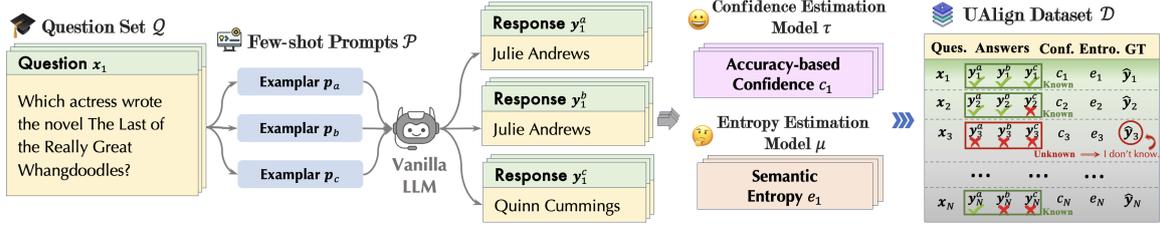


Figure 2: Illustration of UALIGN dataset preparation process.

## 2.1 Dataset Preparation

### 2.1.1 Responses Sampling Strategy

As in Fig. 2, to explore the knowledge boundary of the LLM given a question, we sample multiple responses by repeating the generation procedure several times. In this phase, the preparation process can be represented in a tuple  $(\mathcal{Q}, \mathcal{P}, \mathcal{A})$ .  $\mathcal{Q}$  contains a batch of  $N$  QA pairs  $\{(\mathbf{x}_i, \hat{\mathbf{y}}_i)\}_{i=1}^N$  where  $\mathbf{x}_i$  and  $\hat{\mathbf{y}}_i$  denote the  $i$ -th question and ground-truth answer respectively. To mitigate context sensitivity, we utilize different few-shot prompts in  $\mathcal{P}$  with temperature  $T = 0.2$  to make a trade-off between the accuracy and diversity to represent knowledge boundaries (Gekhman et al., 2024). The few-shot prompt set  $\mathcal{P}$  consists of  $K$  different 1-shot exemplars in this work which is enough for LLMs to generate answers in the correct format. We present the few-shot prompts for sampling on TriviaQA and SciQ datasets as exemplified in Appendix I.

In the  $k$ -th sampling process for the  $i$ -th question  $\mathbf{x}_i$ , we employ each few-shot exemplar  $\mathbf{p}_k \in \mathcal{P}$  with the question  $\mathbf{x}_i$  to the LLM to generate the  $k$ -th response  $\mathbf{y}_i^{(k)}$ . By taking  $K$  times of the sampling process, we can obtain an answer set  $\mathbf{Y}_i = \{\mathbf{y}_i^{(k)}\}_{k=1}^K$  to  $\mathbf{x}_i$ . We set the labels  $\mathbf{Z}_i = \{z_i^{(k)}\}_{k=1}^K$  by comparing each generated answer  $\mathbf{y}_i^{(k)}$  with the ground-truth  $\hat{\mathbf{y}}_i$  to indicate the correctness ( $z_i^{(k)} \in \{0, 1\}$ , 1 for *True* and 0 for *False*). We collect and format the data in  $(\mathbf{x}_i, \mathbf{Y}_i, \mathbf{Z}_i, \hat{\mathbf{y}}_i)$  in an extended dataset and calculate the uncertainty measures subsequently. Note that since fine-tuning LLMs on unknown knowledge will encourage hallucinations (Zhang et al., 2024a; Gekhman et al., 2024), we revise the ground-truth answer to the question with  $z_i^{(k)} = 0, \forall z_i^{(k)} \in \mathbf{Z}_i$  to “Sorry, I don’t know.” to teach LLMs to refuse the questions beyond their knowledge (Zhang et al., 2024a).

### 2.1.2 Uncertainty Measures

In order to quantify the knowledge boundaries, we can leverage some uncertainty estimation methods. The knowledge boundary of LLMs in this work is

defined in two aspects. The first involves the prior judgment to a question  $\mathbf{x}_i$  regardless of the answers (Ren et al., 2023) which indicates the certainty level of  $\mathbf{x}_i$ . The second entails the dispersion measure to the distribution of the generated responses in  $\mathbf{Y}_i$  to  $\mathbf{x}_i$ . Accordingly, we adopt accuracy-based confidence (Xiong et al., 2024) and semantic entropy (Kuhn et al., 2023) to jointly determine and represent the actual knowledge boundary information.

**Accuracy-based Confidence** A natural idea of aggregating varied responses is to measure the accuracy among the candidate outputs to denote confidence scores (Manakul et al., 2023; Xiong et al., 2024). Given a question  $\mathbf{x}_i$ , the accuracy of candidate responses in  $\mathbf{Y}_i$  by comparing with the ground-truth answer  $\hat{\mathbf{y}}_i$  serves as the confidence score  $c_i$ , computed as follows.

$$c_i = \text{Conf}(\mathbf{x}_i) = \frac{1}{K} \sum_{k=1}^K \mathbb{1}(\hat{\mathbf{y}}_i = \mathbf{y}_i^{(k)}) \quad (1)$$

**Semantic Entropy** Due to the variable length and semantically equivalent generated sequences in sentence-level output spaces, Kuhn et al. (2023) proposes semantic entropy to capture uncertainty on the semantic level to quantify the degree of dispersion of sentence meanings. The semantic entropy  $e_i$  given  $\mathbf{x}_i$  and  $\mathbf{Y}_i$  is calculated as

$$p(s|\mathbf{x}_i) = \frac{1}{K} \sum_{k=1}^K \mathbb{1}[\mathbf{y}_i^{(k)} \in s] \quad (2)$$

$$e_i = \text{SE}(\mathbf{x}_i) = - \sum_s p(s|\mathbf{x}_i) \log p(s|\mathbf{x}_i) \quad (3)$$

where  $s$  denotes a set of sentences in semantic equivalent space. As illustrated in Fig. 1, semantic entropy is calculated by clustering semantically equivalent responses, as a measure to quantify the dispersion of generations to confirm the correct answer despite the low confidence, which will be further analyzed with the experimental results in Sec. 4.2. We calculate the confidence score and semantic entropy for both known and unknown questions.

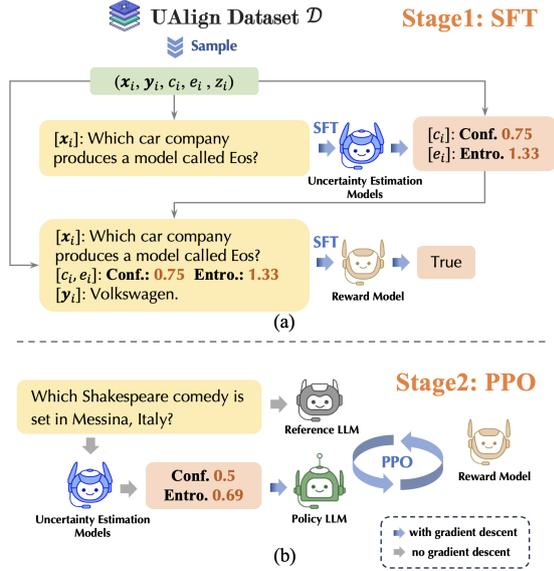


Figure 3: Illustration of (a) SFT and (b) PPO alignment processes of UALIGN framework. Note that for simplicity, we only present one estimation model in the figure but there are actually two.

### Algorithm 1 UALIGN Training Algorithm

- 1: **Input:** UALIGN dataset  $\mathcal{D}$ , uncertainty models  $\tau, \mu$ , reward model  $\theta$ , initial policy  $\pi_o$ .
- 2: **Output:** Optimized policy  $\pi_\theta$ .
- 3: **Stage 1: UALIGN SFT**
- 4: Train uncertainty models  $\tau, \mu$  on  $\mathcal{D}$  to predict  $c_i, e_i$  by feeding  $\mathbf{x}_i$  using Eq. 4 and 5.
- 5: Train reward model  $\theta$  on  $\mathcal{D}$  to predict  $z_i$  by feeding  $\mathbf{x}_i, c_i, e_i, \mathbf{y}_i^{(k)}$  using Eq. 6.
- 6: **Stage 2: UALIGN PPO**
- 7: Collect reward  $r$  including the reward signal  $r_1$  by  $\theta$  and KL-penalty  $r_2$  between policy  $\pi_\theta$  and initial policy  $\pi_o$  as Eq. 7.
- 8: Update policy  $\pi_\theta$  using the collected reward  $r$ .

Then we update a UALIGN dataset  $\mathcal{D}$  by formatting the  $i$ -th sample in  $(\mathbf{x}_i, \mathbf{Y}_i, \mathbf{Z}_i, \hat{\mathbf{y}}_i, c_i, e_i)$ .

## 2.2 UALIGN Training Process

### 2.2.1 UALIGN SFT: Uncertainty Estimation and Reward Models Training

As presented in Fig. 3 (a) and Algorithm 1, given dataset  $\mathcal{D}$ , UALIGN SFT is to train uncertainty estimation models to explicitly learn the two estimations given specific questions. Uncertainty estimation models of  $\tau$  and  $\mu$  are utilized to predict the confidence score and semantic entropy respectively, which are continuously used to train a reward model. When training  $\tau$  and  $\mu$ , we only feed a question  $\mathbf{x}_i$  to the models to generate two uncer-

tainty estimations. The training objectives are to minimize the cross-entropy losses  $\mathcal{L}_\tau$  and  $\mathcal{L}_\mu$  as

$$\arg \min_{\tau} \mathcal{L}_\tau, \arg \min_{\mu} \mathcal{L}_\mu, \quad (4)$$

$$\mathcal{L}_\tau = -\mathbb{E}_{(\mathbf{x}_i, c_i) \sim \mathcal{D}} [\log p_\tau(c_i | \mathbf{x}_i)] \quad (4)$$

$$\mathcal{L}_\mu = -\mathbb{E}_{(\mathbf{x}_i, e_i) \sim \mathcal{D}} [\log p_\mu(e_i | \mathbf{x}_i)] \quad (5)$$

where the models can explicitly learn and express the uncertainty estimations which represent more accurate knowledge boundary information.

Subsequently, the reward model is introduced as a binary evaluator to determine if a generated answer  $\mathbf{y}_i^{(k)} \in \mathbf{Y}_i$  is correctly conditioned on the question  $\mathbf{x}_i$ , confidence  $c_i$ , and entropy  $e_i$ . Both  $c_i$  and  $e_i$  are explicitly used as additional auxiliary features to improve the accuracy of the reward model. The binary cross-entropy loss  $\mathcal{L}_\theta$  for the reward model  $\theta$  is minimized as follows.

$$\arg \min_{\theta} \mathcal{L}_\theta, \mathcal{L}_\theta = -\mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i^{(k)}, z_i^{(k)}, c_i, e_i) \sim \mathcal{D}} [\mathcal{L}_\theta^{(i)}] \quad (6)$$

$$\mathcal{L}_\theta^{(i)} = -z_i^{(k)} \log p_\theta(z_i^{(k)} | \mathbf{x}_i, c_i, e_i, \mathbf{y}_i^{(k)}) \quad (6)$$

$$-(1 - z_i^{(k)}) \log(1 - p_\theta(z_i^{(k)} | \mathbf{x}_i, c_i, e_i, \mathbf{y}_i^{(k)})) \quad (6)$$

### 2.2.2 UALIGN PPO: Policy Model Training

The UALIGN PPO is to elicit the LLM’s factual expressions to a question with the uncertainty measures using obtained models. Inspired by the progress of reinforcement learning from human feedback (RLHF) technique (Ouyang et al., 2022; Ziegler et al., 2019), we employ proximal policy optimization (PPO) (Schulman et al., 2017) for LLM optimization with the reward model  $\theta$ . As illustrated in Fig. 3 (b), the LLM to be optimized is used as the policy  $\pi_\theta$ . During this phase, we iteratively feed the question  $\mathbf{x}$ , and the predicted confidence  $c$  and entropy  $e$  to both the policy  $\pi_\theta$  and the reference  $\pi_o$ , and the reward function  $r$  will facilitate reliable expressions of  $\mathbf{y}$  of the policy model  $\pi_\theta$ . The training objective is to maximize the following reward function  $r$  as

$$\arg \max_{\pi_\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, c \sim \tau(\mathbf{x}), e \sim \mu(\mathbf{x}), \mathbf{y} \sim \pi_\theta(\mathbf{x}, c, e)} [r] \quad (7)$$

$$r = \underbrace{\theta(\mathbf{x}, \mathbf{y}, c, e)}_{r_1} - \beta \underbrace{\text{KL}[\pi_\theta(\mathbf{x}, c, e) || \pi_o(\mathbf{x})]}_{r_2} \quad (7)$$

where the reward function  $r$  contains a reward signal  $r_1$  from  $\theta$  and a KL-penalty  $r_2$  to make sure the generated answers  $\mathbf{y}$  by policy  $\pi_\theta$  don’t diverge too much from the original policy  $\pi_o$ . The hyperparameter  $\beta$  is the coefficient of KL-penalty.

### 3 Experimental Setting

#### 3.1 Datasets

The UALIGN training set is comprised of three widely used knowledge-intensive QA datasets: **TriviaQA (TVQA)** (Joshi et al., 2017) which contains closed-book trivia QA pairs to gauge models’ factual knowledge, **SciQ** (Johannes Welbl, 2017) requiring scientific professional knowledge, and **NQ-Open** (Kwiatkowski et al., 2019) which is constructed by Google Search queries along with annotated short answers or documents.

For testing, we evaluate the in-domain (ID) performance on the corresponding validation/test sets and generalization on an out-of-domain (OOD) test set **LSQA** (Xue et al., 2024) which contains multilingual language-specific QA pairs. More dataset details and statistics are presented in Appendix B.

#### 3.2 Evaluation Metrics

To evaluate the reliability of LLMs, we employ two metrics: *Precision (Prec.)* and *Truthfulness (Truth.)*. *Precision* is defined as the proportion of correctly answered questions among all the known questions, representing LLMs’ ability to accurately express their known factual knowledge. *Truthfulness* represents the proportion of the sum of correctly answered known questions and refused unknown questions among all questions, indicating the honesty level of the LLMs.

To ascertain the correctness of the LLM-generated answer  $y$  with the ground truth  $\hat{y}$ , we employ a string-matching approach. Exact matching (EM) of  $y \equiv \hat{y}$  always misjudges some correct answers with slight distinctions on such closed-book QA tasks. Therefore, we replace EM with a variant of  $y \in \hat{y} \vee \hat{y} \in y$  to evaluate the accuracy. The specific illustrations of evaluation formulas and comparisons of several EM variants we tested with human evaluations are in Appendix C.

#### 3.3 Baselines

We present several baselines in four categories below. To clearly delineate the differences between our proposed method and other baselines, we have illustrated all methods in Fig. 7 in Appendix D.

**Prompt-based** We present two prompt-based baselines namely In-Context Learning (**ICL**), In-Context Learning with Refusal Examples (**ICL-IDK**), and In-Context Learning Chain-of-Thought (**ICL-CoT**) (Wei et al., 2022). The few-shot prompt templates are presented in Appendix E.

**SFT-based** We employ standard Supervised Fine-Tuning (**SFT**) by training an LLM to generate answers for all questions. We also introduce **R-Tuning** (Zhang et al., 2024a) which teaches LLM to refuse their unknown questions.

**RL-based** Following RLHF technique (Ouyang et al., 2022), we first train a reward model to determine correctness by SFT. Then we employ PPO to optimize the policy model with the reward model (**RL-PPO**). We also introduce an advanced variant called reinforcement learning from knowledge feedback (**RLKF**) (Liang et al., 2024) which leverages knowledge probing and consistency checking to train the reward model. Following Zhang et al. (2024b); Tian et al. (2024); Lin et al. (2024), we also construct the factuality preference dataset to conduct direct preference optimization (**RL-DPO**) to enhance the factuality of LLMs.

**Inference-based** Another branch of work focuses on shifting the output distribution to improve factuality during inference. Li et al. (2023) (**ITI**) intervenes in the activations in attention heads to the “truthfulness” direction.

#### 3.4 Implementation Details

Experiments are conducted on two LLMs: **Llama-3-8B** (Llama-3)<sup>2</sup> (AI@Meta, 2024) and **Mistral-7B** (Mistral)<sup>3</sup> (Jiang et al., 2023). When preparing the UALIGN dataset, we sample 10 responses for each question on  $K = 10$  different 1-shot prompts. The sampling temperature  $T$  is set to 0.2 to achieve a trade-off between the diversity and factuality of the answer set. During training, all the LLMs are trained using LoRA (Hu et al., 2022) with rank  $r = 16$ . Both the uncertainty estimation models and the reward model utilize the vanilla LLM as their bases and are trained using LoRA with rank  $r = 4$ . ADAM parameter update is used in a mini-batch mode. Uncertainty estimation models and the reward model are trained using SFT on the UALIGN dataset. The UALIGN PPO algorithm and all the RL-based baselines are implemented by trl<sup>4</sup>. All training hyper-parameters are presented in Appendix F. When decoding, the temperature is also set to 0.2 to be consistent with the sampling setting. All the experiments are conducted on  $4 \times$  NVIDIA A100-40GB GPUs.

<sup>2</sup><https://huggingface.co/meta-llama/Meta-Llama-3-8B>

<sup>3</sup><https://huggingface.co/mistralai/Mistral-7B-v0.1>

<sup>4</sup><https://github.com/huggingface/trl>

Method	TVQA (ID)		SciQ (ID)		NQ-Open (ID)		Avg. (ID)		LSQA (OOD)	
	Prec. $\uparrow$	Truth. $\uparrow$								
<b>Llama-3-8B</b>										
ICL	76.15	56.55	70.43	44.30	50.28	20.11	65.62	40.32	77.35	52.98
ICL-IDK	69.17	54.10	68.36	43.00	45.43	20.72	60.98	39.27	66.67	50.24
ICL-CoT	66.68	53.37	72.34	45.90	<b>57.34</b>	23.60	65.45	40.95	73.96	49.37
SFT	70.80	52.57	72.18	45.40	41.41	16.57	61.46	38.18	68.09	46.63
R-Tuning	72.93	55.44	71.38	44.90	47.81	18.12	64.04	39.48	71.54	52.15
RL-PPO	76.32	55.19	75.70	45.80	54.07	24.19	68.03	41.72	72.18	48.43
RL-DPO	72.08	53.96	71.23	44.20	49.65	19.18	64.32	39.11	71.09	48.88
RLKF	77.12	56.07	72.36	44.90	54.86	22.15	68.11	41.04	74.95	52.46
ITI	71.09	53.97	72.35	43.80	43.20	17.13	62.21	38.30	68.52	46.99
UALIGN	<b>79.14</b>	57.04	<b>76.44</b>	<b>48.00</b>	56.60	26.09	<b>70.72</b>	43.71	<b>79.56</b>	<b>55.88</b>
(w/o Conf.)	74.13	54.45	74.05	45.00	54.19	23.60	67.45	41.01	74.25	52.06
(w/o Entro.)	78.43	<b>57.69</b>	75.39	47.50	56.68	<b>27.56</b>	70.16	<b>44.25</b>	76.14	54.43
<b>Mistral-7B</b>										
ICL	77.92	55.14	68.62	42.20	52.09	17.95	66.21	38.43	74.09	47.71
ICL-IDK	72.59	51.37	63.74	39.20	51.13	17.67	62.48	36.20	72.27	47.32
ICL-CoT	76.73	54.78	71.87	44.20	<b>54.47</b>	18.22	67.69	39.06	<b>79.24</b>	52.59
SFT	74.57	54.77	65.85	42.50	50.82	14.42	63.74	37.08	68.33	44.00
R-Tuning	67.70	52.25	64.44	40.10	46.33	15.52	59.49	36.29	64.67	44.05
RL-PPO	79.23	55.08	71.35	44.10	53.76	19.19	68.11	39.45	74.49	49.67
RL-DPO	72.20	52.98	66.44	41.80	50.95	16.42	63.19	37.06	67.82	43.77
RLKF	80.43	56.92	70.66	43.90	52.09	18.24	67.72	39.68	74.19	49.23
ITI	74.65	55.16	66.90	44.90	51.12	16.68	64.22	38.91	67.73	46.20
UALIGN	<b>82.10</b>	<b>59.05</b>	<b>73.21</b>	<b>46.70</b>	54.17	<b>19.64</b>	<b>70.82</b>	<b>41.79</b>	76.29	<b>52.89</b>
(w/o Conf.)	76.44	55.13	69.84	43.50	50.30	17.88	65.52	38.83	73.15	47.06
(w/o Entro.)	80.18	57.64	72.90	45.60	52.21	18.44	68.43	40.56	75.34	50.15

Table 1: Experiments of Precision (*Prec.*) and Truthfulness (*Truth.*) on four datasets on Llama-3 and Mistral.

## 4 Results and Analysis

### 4.1 Main Experimental Results

We present the results of UALIGN and several baselines on three ID and one OOD test sets as shown in Table 1. Several findings are listed below.

**Reliability** Significant improvements are consistently achieved on diverse datasets using the proposed UALIGN framework over other baseline methods on both Llama-3 and Mistral. We highlight the supreme Precision and Truthfulness performance using grey highlights among the all baselines of each column in Table 1. The core idea of our UALIGN framework is the utilization of uncertainty estimation models. Compared with the most relevant baselines of RL-PPO and RLKF, both the reward model and policy model in UALIGN generate predictions and responses conditioned on uncertainty estimations regarding the knowledge boundaries to questions, thereby yielding better reliability performance. It can be attributed that by explicitly appending uncertainty measures following the question, LLMs can assist LLMs in eliciting more accurate responses based on intrinsic knowledge boundary representations.

**Generalization** We also introduced an OOD test set to assess the generalization capability of the

Conf.	Entro.	TVQA	ID SciQ	NQ-Open	OOD LSQA
<b>Llama-3-8B</b>					
$\times$	$\times$	82.31	79.00	67.45	70.12
$\checkmark$	$\times$	85.41	84.30	70.37	<b>75.09</b>
$\times$	$\checkmark$	82.05	77.90	67.85	70.40
$\checkmark$	$\checkmark$	<b>86.73</b>	<b>86.40</b>	<b>72.00</b>	74.59
<b>Mistral-7B</b>					
$\times$	$\times$	84.53	77.30	65.24	68.31
$\checkmark$	$\times$	86.80	79.50	72.10	72.95
$\times$	$\checkmark$	85.24	74.60	66.64	71.22
$\checkmark$	$\checkmark$	<b>88.06</b>	<b>79.80</b>	<b>75.14</b>	<b>73.61</b>

Table 2: Accuracy of reward model varying different uses of uncertainty measures Conf. and Entro. in UALIGN dataset on Llama-3 and Mistral.

UALIGN method. The results in Table 1 indicate that most training-based baselines (SFT, RL, Inference) are unstable and result in performance decreasing compared with prompt-based baselines when generalizing on the OOD test set. However, comparable reliability performances are obtained on two LLMs using the proposed UALIGN in comparison with prompt-based methods, demonstrating strong generalization capability.

### 4.2 Effects of Uncertainty Estimation Models

**Setting** To investigate the effects of introducing uncertainty estimations as input features to reward models, we report the accuracy of reward models

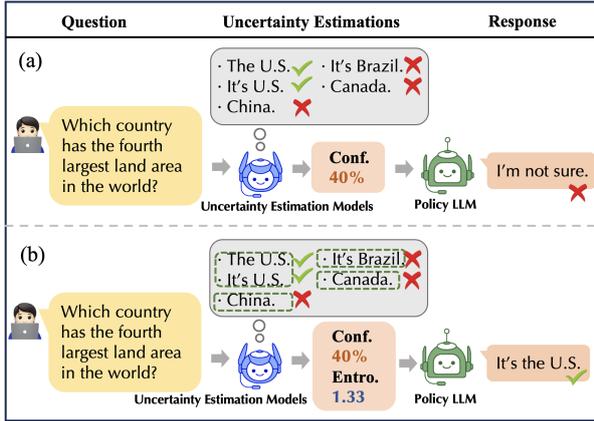


Figure 4: Illustration of the effects of different uses of uncertainty estimations under varying knowledge boundaries perceived by LLMs.

that vary in different uses of two measures on ID and OOD tasks. The reward models are trained on the UALIGN dataset on both Llama-3 and Mistral.

**Results** As in Table 2, we present the results of the accuracy of reward models. Significant accuracy improvements of reward models are obtained that predominantly benefit from the use of confidence scores across both ID and OOD test sets on two LLMs, validating the effectiveness of our proposed UALIGN framework. The isolated use of semantic entropy does not guarantee a stable improvement but may even lead to a performance decrease on some test sets. However, when semantic entropy is employed in combination with confidence measures, it can facilitate further enhancements, achieving optimal results across most test sets as highlighted grey cells for two LLMs.

**Analysis** In the UALIGN framework, both confidence score and semantic entropy are introduced to quantify the intrinsic knowledge boundary of LLMs to questions. The explicit introduction of the knowledge boundary representations in prompts can be regarded as the added thinking step like CoT. The combined use of confidence and semantic entropy can achieve supreme prediction performance in Table 2. We illustrate the mechanism as follows.

As demonstrated in Fig. 4 (a), by sampling multiple responses to a question, we can approximate LLM’s intrinsic knowledge boundary, where the certainty level of the answer “The U.S.” is 40%. In previous work (Zhang et al., 2024a) which only considers the confidence level, the correct answer that the LLM knows but is not sure will be discarded and the LLM will refuse to answer. How-

ever, as in Fig. 4 (b), the LLM can perceive that even though its certainty level to the correct answer is low, other answers are more uncertain and the dispersion level of answers is relatively high which is quantified by semantic entropy. After UALIGN PPO training, the ability to generate correct answers conditioned on questions and estimations is well enhanced. As a result, the correct but unsure knowledge will be elicited in the responses.

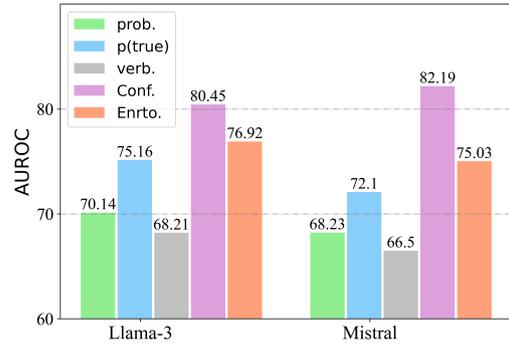


Figure 5: Results of AUROC $\uparrow$  of several uncertainty estimation methods on TVQA using Llama-3 and Mistral.

### 4.3 Reliability of Uncertainty Estimations

**Setting** Evaluating the performance of confidence score and semantic entropy is essential to the UALIGN method. We present the AUROC (Detailed in Appendix C) results of two estimations in comparison with three confidence/uncertainty estimation methods (one probability-based method (Prob.), two prompt-based methods including p(True) and verbalized (Verb.) as illustrated in Fig. 8) on TriviaQA on two LLMs. Results on other datasets are remained in Appendix H. Details of baseline estimation baselines are presented in Sec. 5, Appendix G, and Fig. 8.

**Results** In Fig. 5, both the confidence and entropy prediction consistently outperform other baseline uncertainty estimation methods. Optimal AUROC performances are obtained using confidence on both Llama-3 (80.45) and Mistral (82.19).

**Analysis** After UALIGN SFT stage, the uncertainty estimation models are converged on the UALIGN dataset to predict both confidence and entropy, indicating the models possess the ability to predict the two measures. Practically, our utilized confidence and semantic entropy incorporate the advantages of both sampling- and training-based uncertainty estimations. Multiple sampling can better approximate the actual knowledge boundaries

of LLMs, while the training-based approach enables the LLMs to learn to perceive their intrinsic knowledge boundaries. Compared to other baselines that suffer from overconfidence issues with low AUROC scores, our utilized methods yield more reliable estimates, thereby ensuring improved performance for both the reward model and the policy model in the following stages.

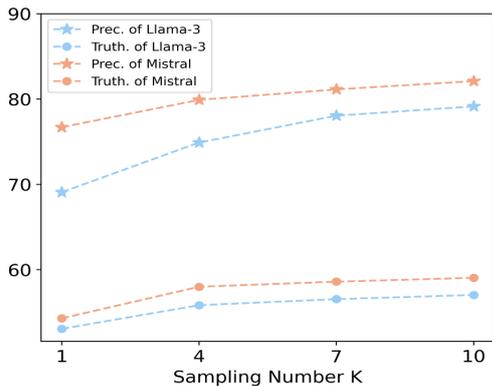


Figure 6: Experiments of Precision (*Prec.*) and Truthfulness (*Truth.*) of various sampling number  $K$  of 1, 4, 7, and 10 on TVQA on Llama-3 and Mistral.

#### 4.4 Effects of Sampling Number

**Setting** The sampling number  $K$  is a crucial hyper-parameter in the UALIGN method. Different values of  $K$  can significantly affect the precision of the knowledge boundary measurements. To evaluate the effects, we compare performances using various  $K$  of 1, 4, 7, and 10. Experimental results on TVQA are presented in Fig. 6 in Appendix H.

**Findings** The experiments indicate that when using small sampling numbers, increasing the  $K$  leads to significant improvements in both precision and truthfulness. However, as  $K$  continues to increase, the reliability improvement tends to plateau, exhibiting convergence. Since further increasing  $K$  requires substantial computational costs, we discard conducting experiments with larger  $K$ .

**Analysis** The results in Fig. 6 demonstrate that while the sampling number  $K$  increases linearly, the performance improvements are non-linear. This may be attributed to utilizing non-linear metrics, or it could suggest that  $K = 10$  can approximate the actual knowledge boundaries, resulting in a gradual slowdown in performance gains. Consequently, setting  $K$  to 10 in this work makes a trade-off between performance gains and computation expense.

## 5 Related Works

**Knowledge Boundary** Previous works investigate the knowledge boundary to identify the known level of a knowledge piece of LLMs by quantifying the confidence or uncertainty estimations like output consistency (Cheng et al., 2024), prompting methods (Ren et al., 2023) or knowledge probing (Ji et al., 2024). Generally, knowledge boundary measures derive from uncertainty estimations.

**Uncertainty Estimation for LLMs** We categorize uncertainty estimation methods on LLMs into four classes as illustrated in Figure 8. ① *Likelihood-based methods* Vazhentsev et al. (2023) directly quantify sentence uncertainty over token probabilities; ② *Prompting-based methods* instruct LLMs to express uncertainty in words (Lin et al., 2022a; Xiong et al., 2024) or to self-evaluate its correctness on  $p(\text{True})$  (Kadavath et al., 2022); ③ *Sampling-based methods* aggregate sampled responses to calculate consistency (Xiong et al., 2024) or semantic entropy (Kuhn et al., 2023); ④ *Training-based methods* (Lin et al., 2022a) propose to train LLMs to improve linguistic uncertainty expressions.

**Factuality Alignment** LLM alignment is to guide human preference through Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022; Bai et al., 2022a). Distinct from recent studies that apply RL to improve LLMs’ factuality (Zhang et al., 2024b; Lin et al., 2024; Liang et al., 2024; Xu et al., 2024a), this work improves LLMs’ reliability by explicitly leveraging the uncertainty estimations for LLM alignment.

Due to the space limitation, detailed investigations of related works are shown in Appendix G.

## 6 Conclusion

In this paper, we present a UALIGN framework to explicitly leverage uncertainty estimations to elicit LLMs to accurately express factual knowledge that LLMs cannot constantly answer correctly due to ambiguous knowledge boundaries. We introduce the dataset preparation process and UALIGN training strategies of factuality alignment by incorporating uncertainty estimations of the confidence score and semantic entropy as input features into prompts. Experiments on several knowledge QA tasks affirm the efficacy of UALIGN to enhance the LLMs’ reliability and generalizability, demonstrating significant improvements over various baselines.

## 585 Limitations

586 The limitations and future work of this study are  
587 listed as follows:

588 **Computational Resources:** The current method  
589 for constructing the UALIGN dataset relies on  
590 multiple samplings, requiring substantial computa-  
591 tional cost which linearly increases with the num-  
592 ber of sampling instances  $K$  and significantly con-  
593 strains the scalability expansion of the dataset in  
594 this work. To accurately approximate the knowl-  
595 edge distributions, a higher number of samplings  
596 is typically more beneficial. Here, we may need  
597 to introduce some prior knowledge distributions to  
598 alleviate the computational resource requirements  
599 and reduce the costs.

600 **Task Expansion:** The dataset used in this pa-  
601 per is solely based on factual knowledge QA  
602 tasks, with a simple and fixed template and re-  
603 sponse format. However, the UALIGN methodol-  
604 ogy has not been further validated on other fac-  
605 tual knowledge-based tasks such as open-form  
606 instruction-following tasks, long-form generation  
607 like biography, or even knowledge reasoning tasks,  
608 where the uncertainty estimations remain chal-  
609 lenging. In future works, we plan to extend the  
610 UALIGN framework to open-ended generation  
611 tasks to enhance the LLMs’ factual expressions.

## 612 Acknowledgments

## 613 References

614 Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana  
615 Rezazadegan, Li Liu, Mohammad Ghavamzadeh,  
616 Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Ra-  
617 jendra Acharya, et al. 2021. A review of uncertainty  
618 quantification in deep learning: Techniques, appli-  
619 cations and challenges. *Information fusion*, 76:243–  
620 297.

621 AI@Meta. 2024. [Llama 3 model card](#). *AI@Meta*.

622 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda  
623 Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
624 Stanislav Fort, Deep Ganguli, Tom Henighan,  
625 Nicholas Joseph, Saurav Kadavath, Jackson Kernion,  
626 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac  
627 Hatfield-Dodds, Danny Hernandez, Tristan Hume,  
628 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel  
629 Nanda, Catherine Olsson, Dario Amodei, Tom  
630 Brown, Jack Clark, Sam McCandlish, Chris Olah,  
631 Ben Mann, and Jared Kaplan. 2022a. [Training  
632 a helpful and harmless assistant with reinforce-  
633 ment learning from human feedback](#). *Preprint*,  
634 arXiv:2204.05862.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, 635  
Amanda Askell, Jackson Kernion, Andy Jones, 636  
Anna Chen, Anna Goldie, Azalia Mirhoseini, 637  
Cameron McKinnon, et al. 2022b. Constitutional 638  
ai: Harmlessness from ai feedback. *arXiv preprint* 639  
*arXiv:2212.08073*. 640

Tom Brown, Benjamin Mann, Nick Ryder, Melanie 641  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 642  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 643  
Askell, Sandhini Agarwal, Ariel Herbert-Voss, 644  
Gretchen Krueger, Tom Henighan, Rewon Child, 645  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens 646  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma- 647  
teusz Litwin, Scott Gray, Benjamin Chess, Jack 648  
Clark, Christopher Berner, Sam McCandlish, Alec 649  
Radford, Ilya Sutskever, and Dario Amodei. 2020. 650  
[Language models are few-shot learners](#). In *Ad-  
651 vances in Neural Information Processing Systems*,  
652 volume 33, pages 1877–1901. Curran Associates, 653  
Inc. 654

Jiuhai Chen and Jonas Mueller. 2023. Quantifying un- 655  
certainty in answers from any language model via 656  
intrinsic and extrinsic confidence assessment. *arXiv* 657  
*preprint arXiv:2308.16175*. 658

Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wen- 659  
wei Zhang, Zhangyue Yin, Shimin Li, Linyang Li, 660  
Zhengfu He, Kai Chen, and Xipeng Qiu. 2024. [Can  
661 ai assistants know what they don’t know?](#) *Preprint*,  
662 arXiv:2401.13275. 663

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, 664  
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan 665  
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion 666  
Stoica, and Eric P. Xing. 2023. [Vicuna: An open-  
667 source chatbot impressing gpt-4 with 90%\\* chatgpt  
668 quality](#). 669

Angelos Filos, Sebastian Farquhar, Aidan N Gomez, 670  
Tim GJ Rudner, Zachary Kenton, Lewis Smith, Mil- 671  
lad Alizadeh, Arnoud de Kroon, and Yarin Gal. 672  
2019. Benchmarking bayesian deep learning with di- 673  
abetic retinopathy diagnosis. *Preprint at https://arxiv.  
674 org/abs/1912.10481*. 675

Yarin Gal and Zoubin Ghahramani. 2016. [Dropout as  
676 a bayesian approximation: Representing model un-  
677 certainty in deep learning](#). In *Proceedings of The  
678 33rd International Conference on Machine Learn-  
679 ing*, volume 48 of *Proceedings of Machine Learning  
680 Research*, pages 1050–1059, New York, New York,  
681 USA. PMLR. 682

Yarin Gal et al. 2016. Uncertainty in deep learning. 683  
*Ph.D. Thesis*. 684

Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, 685  
Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. 686  
[Does fine-tuning LLMs on new knowledge encour-  
687 age hallucinations?](#) In *Proceedings of the 2024 Con-  
688 ference on Empirical Methods in Natural Language  
689 Processing*, pages 7765–7784, Miami, Florida, USA.  
690 Association for Computational Linguistics. 691

692	Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. 2023. <a href="#">A survey of language model confidence estimation and calibration</a> . <i>Preprint</i> , arXiv:2311.08298.	747
693		748
694		749
695		750
696	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. <a href="#">On calibration of modern neural networks</a> . <i>Preprint</i> , arXiv:1706.04599.	751
697		752
698		753
699	Haixia Han, Tingyun Li, Shisong Chen, Jie Shi, Chengyu Du, Yanghua Xiao, Jiaqing Liang, and Xin Lin. 2024. <a href="#">Enhancing confidence expression in large language models through learning from past experience</a> . <i>Preprint</i> , arXiv:2404.10315.	754
700		755
701		756
702		757
703		758
704	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. <a href="#">LoRA: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .	759
705		760
706		761
707		762
708		763
709	Ziwei Ji, Delong Chen, Etsuko Ishii, Samuel Cahyawijaya, Yejin Bang, Bryan Wilie, and Pascale Fung. 2024. <a href="#">LLM internal states reveal hallucination risk faced with a query</a> . In <i>Proceedings of the 7th Black-boxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP</i> , pages 88–104, Miami, Florida, US. Association for Computational Linguistics.	764
710		765
711		766
712		767
713		768
714		769
715		770
716	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. <a href="#">Mistral 7b</a> . <i>Preprint</i> , arXiv:2310.06825.	771
717		772
718		773
719		774
720		775
721		776
722		777
723		778
724	Matt Gardner Johannes Welbl, Nelson F. Liu. 2017. <a href="#">Crowdsourcing multiple choice science questions</a> . In <i>arXiv</i> .	779
725		780
726		781
727	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. <a href="#">TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension</a> . 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.	782
728		783
729		784
730		785
731		786
732		787
733		788
734	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. <a href="#">Language models (mostly) know what they know</a> . <i>arXiv preprint arXiv:2207.05221</i> .	789
735		790
736		791
737		792
738		793
739		794
740	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. <a href="#">Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	795
741		796
742		797
743		798
744		799
745	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti,	800
746		801
	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: A benchmark for question answering research</a> . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	747
		748
		749
		750
		751
		752
		753
	Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. <a href="#">Simple and scalable predictive uncertainty estimation using deep ensembles</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	754
		755
		756
		757
		758
	Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. 2018. <a href="#">A simple unified framework for detecting out-of-distribution samples and adversarial attacks</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 31. Curran Associates, Inc.	759
		760
		761
		762
		763
	Kenneth Li, Oam Patel, Fernanda Vi��gas, Hanspeter Pfister, and Martin Wattenberg. 2023. <a href="#">Inference-time intervention: Eliciting truthful answers from a language model</a> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	764
		765
		766
		767
		768
	Siheng Li, Cheng Yang, Taiqiang Wu, Chufan Shi, Yuji Zhang, Xinyu Zhu, Zesen Cheng, Deng Cai, Mo Yu, Lemao Liu, Jie Zhou, Yujiu Yang, Ngai Wong, Xixin Wu, and Wai Lam. 2024. <a href="#">A survey on the honesty of large language models</a> . <i>Preprint</i> , arXiv:2409.18786.	769
		770
		771
		772
		773
	Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024. <a href="#">Learning to trust your feelings: Leveraging self-awareness in LLMs for hallucination mitigation</a> . In <i>Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP</i> , pages 44–58, Bangkok, Thailand. Association for Computational Linguistics.	774
		775
		776
		777
		778
		779
		780
	Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen tau Yih, and Xilun Chen. 2024. <a href="#">Flame: Factuality-aware alignment for large language models</a> . <i>Preprint</i> , arXiv:2405.01525.	781
		782
		783
		784
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a. <a href="#">Teaching models to express their uncertainty in words</a> . <i>Transactions on Machine Learning Research</i> .	785
		786
		787
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022b. <a href="#">TruthfulQA: Measuring how models mimic human falsehoods</a> . 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.	788
		789
		790
		791
		792
		793
	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. <a href="#">Generating with confidence: Uncertainty quantification for black-box large language models</a> . <i>arXiv preprint arXiv:2305.19187</i> .	794
		795
		796
		797
	Genglin Liu, Xingyao Wang, Lifan Yuan, Yangyi Chen, and Hao Peng. 2024. <a href="#">Examining llms’ uncertainty expression towards questions outside parametric knowledge</a> . <i>Preprint</i> , arXiv:2311.09731.	798
		799
		800
		801

802	Jianqiao Lu, Zhiyang Dou, Hongru Wang, Zeyu Cao,	John Schulman, Filip Wolski, Prafulla Dhariwal,	858
803	Jianbo Dai, Yingjia Wan, and Zhijiang Guo. 2024.	Alec Radford, and Oleg Klimov. 2017. Proxi-	859
804	<a href="#">Autopsv: Automated process-supervised verifier.</a>	mal policy optimization algorithms. <i>arXiv preprint</i>	860
805	<i>Preprint</i> , arXiv:2405.16802.	<i>arXiv:1707.06347</i> .	861
806	Qing Lyu, Kumar Shridhar, Chaitanya Malaviya,	Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg	862
807	Li Zhang, Yanai Elazar, Niket Tandon, Marianna	Durrett. 2023. <a href="#">A long way to go: Investigating length</a>	863
808	Apidianaki, Mrinmaya Sachan, and Chris Callison-	<a href="#">correlations in rlhf.</a> <i>Preprint</i> , arXiv:2310.03716.	864
809	Burch. 2024. <a href="#">Calibrating large language models with</a>		
810	<a href="#">sample consistency.</a> <i>Preprint</i> , arXiv:2402.13904.		
811	Andrey Malinin and Mark Gales. 2021. <a href="#">Uncertainty</a>	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christo-	865
812	<a href="#">estimation in autoregressive structured prediction.</a>	pher D Manning, and Chelsea Finn. 2024. <a href="#">Fine-</a>	866
813	In <i>International Conference on Learning Representa-</i>	<a href="#">tuning language models for factuality.</a> In <i>The Twelfth</i>	867
814	<i>tions.</i>	<i>International Conference on Learning Representa-</i>	868
815	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023.	<i>tions.</i>	869
816	<a href="#">SelfCheckGPT: Zero-resource black-box hallucina-</a>	Katherine Tian, Eric Mitchell, Allan Zhou, Archit	870
817	<a href="#">tion detection for generative large language models.</a>	Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn,	871
818	In <i>Proceedings of the 2023 Conference on Empiri-</i>	and Christopher Manning. 2023a. <a href="#">Just ask for cali-</a>	872
819	<i>cal Methods in Natural Language Processing</i> , pages	<a href="#">bration: Strategies for eliciting calibrated confidence</a>	873
820	9004–9017, Singapore. Association for Computa-	<a href="#">scores from language models fine-tuned with human</a>	874
821	tional Linguistics.	<a href="#">feedback.</a> In <i>Proceedings of the 2023 Conference</i>	875
822	Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-	<i>on Empirical Methods in Natural Language Process-</i>	876
823	Lan Boureau. 2022. <a href="#">Reducing conversational agents’</a>	<i>ing</i> , pages 5433–5442, Singapore. Association for	877
824	<a href="#">overconfidence through linguistic calibration.</a> <i>Trans-</i>	Computational Linguistics.	878
825	<i>actions of the Association for Computational Linguis-</i>	Katherine Tian, Eric Mitchell, Allan Zhou, Archit	879
826	<i>tics</i> , 10:857–872.	Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn,	880
827	Kenton Murray and David Chiang. 2018. <a href="#">Correcting</a>	and Christopher D Manning. 2023b. <a href="#">Just ask for cali-</a>	881
828	<a href="#">length bias in neural machine translation.</a> In <i>Proceed-</i>	<a href="#">bration: Strategies for eliciting calibrated confidence</a>	882
829	<i>ings of the Third Conference on Machine Translation:</i>	<a href="#">scores from language models fine-tuned with human</a>	883
830	<i>Research Papers</i> , pages 212–223, Brussels, Belgium.	<a href="#">feedback.</a> <i>arXiv preprint arXiv:2305.14975</i> .	884
831	Association for Computational Linguistics.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	885
832	Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	886
833	Martinen. 2024. <a href="#">Kernel language entropy: Fine-</a>	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	887
834	<a href="#">grained uncertainty quantification for llms from se-</a>	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	888
835	<a href="#">mantic similarities.</a> <i>Preprint</i> , arXiv:2405.20003.	Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open</a>	889
836	OpenAI. 2023. <a href="#">Gpt-4 technical report.</a> <i>Preprint</i> ,	<a href="#">and efficient foundation language models.</a> <i>Preprint</i> ,	890
837	arXiv:2303.08774.	arXiv:2302.13971.	891
838	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jian-	892
839	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	shu Chen, and Dong Yu. 2023. <a href="#">A stitch in time saves</a>	893
840	Sandhini Agarwal, Katarina Slama, Alex Ray, John	nine: <a href="#">Detecting and mitigating hallucinations of</a>	894
841	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	llms by validating low-confidence generation. <i>arXiv</i>	895
842	Maddie Simens, Amanda Askell, Peter Welinder,	<i>preprint arXiv:2307.03987</i> .	896
843	Paul F Christiano, Jan Leike, and Ryan Lowe. 2022.	Artem Vazhentsev, Akim Tsvigun, Roman Vashurin,	897
844	<a href="#">Training language models to follow instructions with</a>	Sergey Petrakov, Daniil Vasilev, Maxim Panov,	898
845	<a href="#">human feedback.</a> In <i>Advances in Neural Information</i>	Alexander Panchenko, and Artem Shelmanov. 2023.	899
846	<i>Processing Systems</i> , volume 35, pages 27730–27744.	<a href="#">Efficient out-of-domain detection for sequence to se-</a>	900
847	Curran Associates, Inc.	<a href="#">quence models.</a> In <i>Findings of the Association for</i>	901
848	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano	<i>Computational Linguistics: ACL 2023</i> , pages 1430–	902
849	Ermon, Christopher D. Manning, and Chelsea Finn.	1454, Toronto, Canada. Association for Computa-	903
850	2023. <a href="#">Direct preference optimization: Your lan-</a>	tional Linguistics.	904
851	<a href="#">guage model is secretly a reward model.</a> <i>Preprint</i> ,	Hao Wang and Dit-Yan Yeung. 2020. <a href="#">A survey on</a>	905
852	arXiv:2305.18290.	<a href="#">bayesian deep learning.</a> <i>ACM Comput. Surv.</i> , 53(5).	906
853	Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin	Hongru Wang, Boyang Xue, Baohang Zhou, Tianhua	907
854	Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen,	Zhang, Cunxiang Wang, Guanhua Chen, Huimin	908
855	and Haifeng Wang. 2023. <a href="#">Investigating the factual</a>	Wang, and Kam fai Wong. 2024. <a href="#">Self-dc: When to re-</a>	909
856	<a href="#">knowledge boundary of large language models with</a>	<a href="#">trieve and when to generate? self divide-and-conquer</a>	910
857	<a href="#">retrieval augmentation.</a> <i>Preprint</i> , arXiv:2307.11019.	<a href="#">for compositional unknown questions.</a> <i>Preprint</i> ,	911
		arXiv:2402.13514.	912

913	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. <a href="#">Chain of thought prompting elicits reasoning in large language models</a> . In <i>Advances in Neural Information Processing Systems</i> .	Hanning Zhang, Shizhe Diao, Yong Lin, Yi R. Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. <a href="#">R-tuning: Instructing large language models to say ‘i don’t know’</a> . <i>Preprint</i> , arXiv:2311.09677.	968 969 970 971 972
918	Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2022. <a href="#">Uncertainty quantification with pre-trained language models: A large-scale empirical analysis</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 7273–7284, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024b. <a href="#">Self-alignment for factuality: Mitigating hallucinations in llms via self-evaluation</a> . <i>Preprint</i> , arXiv:2402.09267.	973 974 975 976 977
926	Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie Fu, Junxian He, and Bryan Hooi. 2024. <a href="#">Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. 2023. <a href="#">Navigating the grey area: How expressions of uncertainty and overconfidence affect language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5506–5524, Singapore. Association for Computational Linguistics.	978 979 980 981 982 983 984
931	Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai Fan, Lu Chen, and Kai Yu. 2024a. <a href="#">Rejection improves reliability: Training LLMs to refuse unknown questions using RL from knowledge feedback</a> . In <i>First Conference on Language Modeling</i> .	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. <i>arXiv preprint arXiv:1909.08593</i> .	985 986 987 988 989
936	Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024b. <a href="#">Knowledge conflicts for LLMs: A survey</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 8541–8565, Miami, Florida, USA. Association for Computational Linguistics.		
943	Boyang Xue, Shoukang Hu, Junhao Xu, Mengzhe Geng, Xunying Liu, and Helen Meng. 2022. <a href="#">Bayesian neural network language modeling for speech recognition</a> . <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 30:2900–2917.		
948	Boyang Xue, Hongru Wang, Rui Wang, Sheng Wang, Zezhong Wang, Yiming Du, Bin Liang, and Kam-Fai Wong. 2024. <a href="#">A comprehensive study of multilingual confidence estimation on large language models</a> . <i>Preprint</i> , arXiv:2402.13606.		
953	Boyang Xue, Weichao Wang, Hongru Wang, Fei Mi, Rui Wang, Yasheng Wang, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. 2023. <a href="#">Improving factual consistency for knowledge-grounded dialogue systems via knowledge enhancement and alignment</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 7829–7844, Singapore. Association for Computational Linguistics.		
961	Yuchen Yang, Houqiang Li, Yanfeng Wang, and Yu Wang. 2023. <a href="#">Improving the reliability of large language models by leveraging uncertainty-aware in-context learning</a> . <i>Preprint</i> , arXiv:2310.04782.		
965	Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2024. <a href="#">Alignment for honesty</a> . <i>Preprint</i> , arXiv:2312.07000.		

Notation	Description
$\mathcal{Q}$	Dataset containing $n$ Question-Answering pairs. ( $ \mathcal{Q}  = n$ )
$\mathcal{P}$	Set of few-shot exemplars.
$\mathbf{x}_i$	The $i$ -th question sample in $\mathcal{Q}$ .
$\hat{\mathbf{y}}_i$	The $i$ -th ground-truth answer in $\mathcal{Q}$ .
$\mathbf{y}_i^{(k)}$	The $k$ -th sampled response to the $i$ -th question in $\mathcal{Q}$ .
$\mathbf{p}_k$	$k$ -th few-shot exemplar to sample $\mathbf{y}_i^{(k)}$ .
$\mathbf{Y}_i$	Answering set containing $K$ sampled response $\{\mathbf{y}_i^{(k)}\}$ for the $i$ -th question $\mathbf{x}_i$ .
$z_i^{(k)}$	The label of $\mathbf{y}_i^{(k)}$ ( $z_i^{(k)} \in \{0, 1\}$ , 1 for <i>True</i> and 0 for <i>False</i> ).
$\mathbf{Z}_i$	Label set corresponding to $\mathbf{Y}_i$ .
$c_i$	The confidence score for the $i$ -th question $\mathbf{x}_i$ .
$e_i$	The semantic entropy for the $i$ -th question $\mathbf{x}_i$ .
$\mathcal{D}$	Constructed UALIGN training set containing $N$ tuple samples $(\mathbf{x}_i, \mathbf{Y}_i, \mathbf{Z}_i, \hat{\mathbf{y}}_i, c_i, e_i)$ .
$\tau$	Uncertainty estimation model trained to calculate confidence score by feeding $\mathbf{x}$ .
$\mu$	Uncertainty estimation model trained to calculate semantic entropy by feeding $\mathbf{x}$ .
$\theta$	Binary classifier by feeding $(\mathbf{x}, c, e, \mathbf{y})$ as the reward model.
$\mathcal{L}_{\mathcal{M}}$	Training loss functions for three models respectively where $\mathcal{M} \in \{\tau, \mu, \theta\}$ .
$r$	Final reward signal consisted of reward score $r_1$ and KL-penalty $r_2$ .
$\beta$	Coefficient for the KL-penalty $r_2$ .
$\pi_{\theta}$	Policy model to be optimized using $r$ by PPO.
$\pi_o$	Reference model initialized by the original policy.
$T$	Sampling temperatue.
$K$	Number of sampled responses.
$N$	Number of QA pairs.

Table 3: Summarized notations in this work.

## A Protocols

### A.1 Definition of Notations

The definitions of the notations in this work are summarized in Table 3.

### A.2 Terminology Use

- In this work, ‘‘UALIGN’’ in small caps font specifically indicates the proposed framework, which indicates methodology like UALIGN dataset, UALIGN SFT and UALIGN PPO.

## B Dataset Details

**TriviaQA** The TriviaQA dataset (Joshi et al., 2017)<sup>5</sup> is a comprehensive reading comprehension dataset of QA resource consisting of approximately 650,000 question-answer-evidence triples sourced from 95,000 documents on Wikipedia and various other websites. This dataset is distinguished by its complexity and serves as an effective benchmark for evaluating machine comprehension and open-domain QA systems. Unlike standard QA benchmark datasets, where answers are directly retrievable, TriviaQA presents a more rigorous challenge as it requires deeper inference to derive answers.

<sup>5</sup>[https://huggingface.co/datasets/mandarjoshi/trivia\\_qa](https://huggingface.co/datasets/mandarjoshi/trivia_qa)

When constructing the UALIGN dataset, we pre-process and extract 76,523 QA samples from the TriviaQA training set and 9,960 from the development set to contribute to the UALIGN training and in-domain test set respectively. Since approximating the knowledge distribution of a question requires multiple sampling where the computation cost is linearly increasing with the sampling time  $K$ , to simplify the setup and conserve computation resources, we conducted experiments using half of the training data points from the original dataset.

**SciQ** The SciQ dataset (Johannes Welbl, 2017)<sup>6</sup> contains 13,679 crowd-sourced science exam questions about physics, chemistry and biology, among others. The original dataset was divided, with 11,679 samples allocated as the training set and an additional 1,000 samples designated as the validation set. These were subsequently incorporated into our UALIGN training set and in-domain test set, respectively.

**NQ-Open** The NQ-Open dataset is derived from Natural Question (Kwiatkowski et al., 2019)<sup>7</sup>,

<sup>6</sup><https://huggingface.co/datasets/allenai/sciq>

<sup>7</sup>[https://huggingface.co/datasets/google-research-datasets/nq\\_open](https://huggingface.co/datasets/google-research-datasets/nq_open)

which is a QA dataset consisting of real queries issued to the Google search engine. We employ the training and development set of NQ-Open, which contains 87,925 and 3,610 samples respectively, to further enhance the UALIGN training and in-domain test set. Since data construction is highly expensive, we also randomly sample half of the QA pairs from the source training data. We mix the selected training samples to construct the UALIGN dataset, which is further used for U2Align SFT+PPO training.

**LSQA** The LSQA dataset is a multilingual knowledge-intensive QA dataset pertaining to language-dominant knowledge covering specific social, geographical, and cultural language contexts for the UK & US, France, China, Japan, and Thailand respectively. In this study, we only input the QA pairs in English from each LSQA subset which includes 1,025 samples as the out-of-domain test set.

## C Evaluation Details

**Accuracy** For closed-book QA evaluation, we observe that simply applying EM may misjudge the correct answers. We compare several variants of EM as in Table 4 and report their successful judgments on responses of 20 selected samples that are misjudged using EM, where PEM, RRM, and PREM indicate Positive-EM, Recall-EM, and Positive-Recall-EM and the mathematical explanations are presented in Table 4. Upon human discrimination, EMPR exhibits the lowest failure rate and is therefore selected as the evaluation metric for this work.

Variant	Explanation	# Fail
EM	$\mathbf{y} \equiv \hat{\mathbf{y}}$	20
PEM	$\mathbf{y} \in \hat{\mathbf{y}}$	16
REM	$\hat{\mathbf{y}} \in \mathbf{y}$	6
PREM	$\mathbf{y} \in \hat{\mathbf{y}} \vee \hat{\mathbf{y}} \in \mathbf{y}$	2

Table 4: Number of failed judgments by human check for different EM variants.

**Area Under the Receiver Operator Characteristic Curve (AUROC)** AUROC assesses the effectiveness of confidence estimation (Filos et al., 2019) by quantifying how likely a randomly chosen correct answer possesses a higher confidence score than an incorrect one, yielding a score within the

range of [0, 1], implemented by sklearn toolkit<sup>8</sup>. A higher AUROC score implying higher reliability is preferred.

## D Baseline Details

**Prompt-based** For all in-context learning methods, we extract the examples from the respective training set to mitigate the knowledge distribution shift between different datasets. For example, the demonstrated examples in Appendix I are derived from the TriviaQA training set and are specifically used when inferring on the TriviaQA validation set. For LSQA without the training set, we use the same examples as TriviaQA as their knowledge domains largely overlap.

- **ICL:** Few-shot prompts containing  $m$  examples are utilized for answer generation with temperature  $T = 0.2$  where  $m$  is set to 2 as presented in the Template E.
- **ICL-IDK:** Two examples are included in the few-shot prompt while one is selected from the ICL-used example, and another is an unknown question whose answer is revised to “*Sorry, I don’t know.*” as presented in the Template E.
- **ICL-CoT:** We also employ the Chain-of-Thought in few-shot examples by recalling the relevant knowledge piece of LLMs and incorporating it into thinking steps before answering the question as presented in the Template E.
- **SFT:** The standard supervised fine-tuning (SFT) is implemented by minimizing the negative log-likelihood of the ground-truth  $\hat{\mathbf{y}}$  conditioned on input question  $\mathbf{x}$  on model  $\pi$ .

$$\arg \min_{\pi} \mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(\mathbf{x}_i, \hat{\mathbf{y}}_i) \sim \mathcal{D}} [\log p_{\pi}(\hat{\mathbf{y}}|\mathbf{x})] \quad (8)$$

- **R-Tuning:** R-Tuning (Zhang et al., 2024a) is implemented in the same way as SFT which only revises the ground-truth label of unknown questions to the refusal answers. The unknown questions are determined if all the sampled responses in the UALIGN dataset are incorrect.

<sup>8</sup>[https://github.com/scikit-learn/scikit-learn/blob/main/sklearn/metrics/\\_ranking.py](https://github.com/scikit-learn/scikit-learn/blob/main/sklearn/metrics/_ranking.py)

- **RL-PPO:** Following (Ouyang et al., 2022), we develop the RL-PPO by training a reward model using the LLM-generated incorrect responses as negative samples. Then we conduct the PPO (Schulman et al., 2017) algorithm with the obtained reward model. In other word, the RL-PPO baseline is a variant of UALIGN which discards the uncertainty estimations.
- **RLKF:** Following (Liang et al., 2024), we employ the RLKF baseline by training the reward model on the LLMs’ internal states with the knowledge probes and further conduct PPO using the reward model. The knowledge probing setting and implementations are referred to as Liang et al. (2024).
- **RL-DPO:** All Tian et al. (2024); Lin et al. (2024); Zhang et al. (2024b) focus on long-context generation like biography. We still utilize the LLMs’ generated incorrect responses as negative samples to construct the preference data to conduct the DPO (Rafailov et al., 2023) algorithm.
- **ITI:** We replicate (Li et al., 2023) by training a head probe in the attention layer to intervene in the activations to the “truthfulness” direction. To be consistent with the original work, we also train the head on TruthfulQA (Lin et al., 2022b) with our prepared UALIGN dataset to decode in the “truthfulness” direction. Then we further train the LLM using LoRA by SFT to adapt QA tasks. Therefore, the replicated ITI can be regarded as conducting SFT on LLMs with an additional “truthfulness” head.

## E Prompt Template

```

ICL Prompt

You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.

### Question ###: {demo_question_1}
### Answer ###: {demo_answer_1}

### Question ###: {demo_question_2}
### Answer ###: {demo_answer_2}

### Question ###: {input_question}
### Answer ###:

```

```

ICL-IDK Prompt

You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.

### Question ###: {demo_question_1}
### Answer ###: {demo_answer_1}

### Question ###: {demo_question_2}
### Answer ###: {refusal}

### Question ###: {input_question}
### Answer ###:

```

```

ICL-CoT Prompt

You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.

### Question ###: {demo_question_1}
### Recall ###: {knowledge_1}
### Answer ###: {demo_answer_1}

### Question ###: {demo_question_2}
### Recall ###: {knowledge_2}
### Answer ###: {demo_answer_2}

### Question ###: {input_question}
### Answer ###:

```

## F Training Setting Details

To conserve memory overhead and accelerate computation, all the models are quantified using float16 (fp16) to load and save parameters during both the training and inference phases. During the training stage, the batch sizes for the LLM, uncertainty estimation models, and reward models are set at 4, 16, and 16, respectively. The initial learning rate of 1e-4 is utilized with the 0.05 warm-up ratio and 0.01 weight decay of the ADAM optimizer. We set the training epoch to 2 and ensure that all the models can be trained to convergence by increasing additional training steps if necessary. The dropout rate is set at 0.05 during all model updates to reduce overfitting. In the RL phase, all the hyper-parameters related to PPO algorithm are default values by the tr1 PPOConfig recipe<sup>9</sup> except the epoch, learning rate, and batch size which are set at 2, 1e-5, and 2, respectively.

<sup>9</sup>[https://github.com/huggingface/trl/blob/main/trl/trainer/po\\_config.py](https://github.com/huggingface/trl/blob/main/trl/trainer/po_config.py)

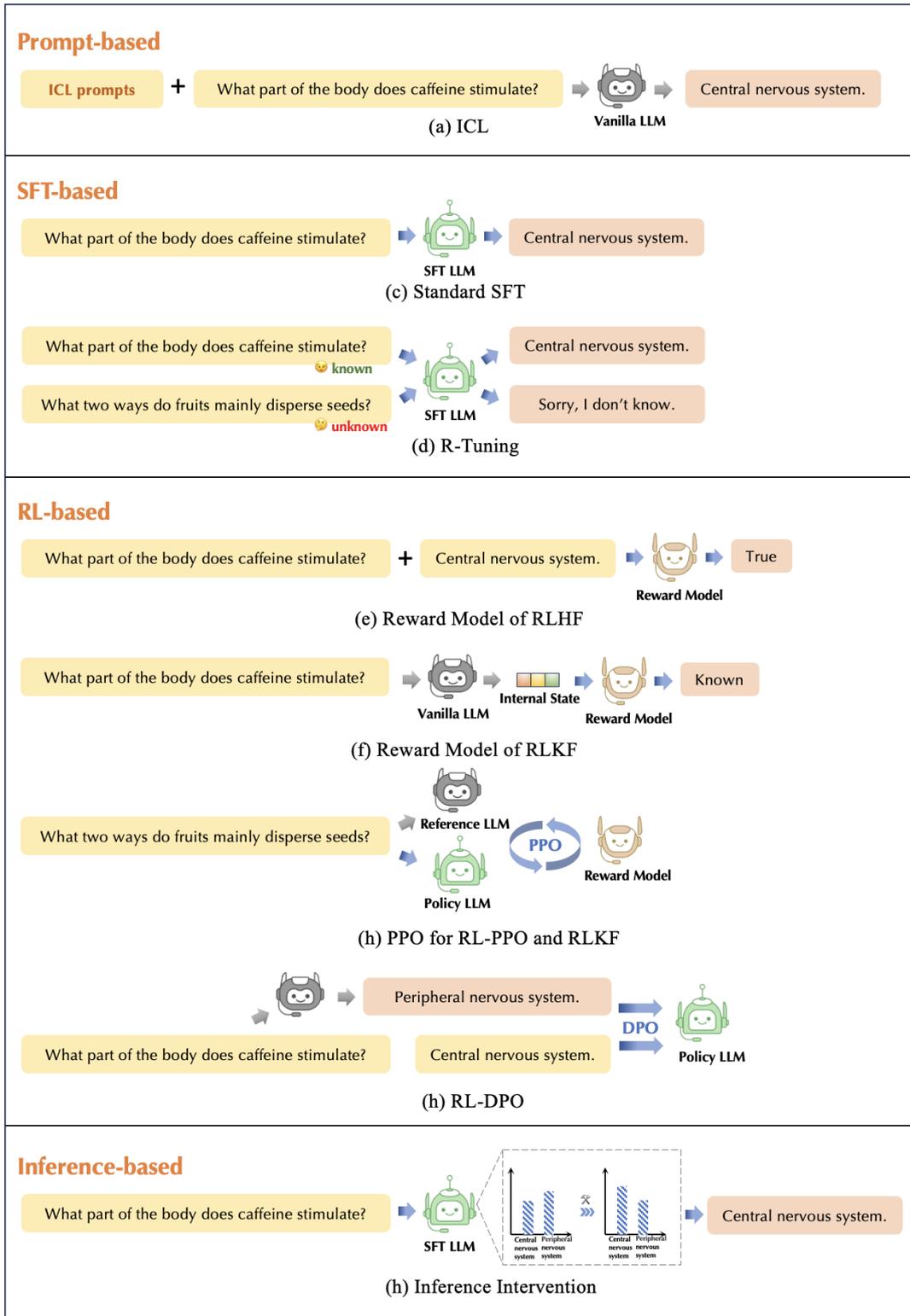


Figure 7: Illustration of several baselines as in Sec. 3.3.

## G Detailed Related Works

### G.1 Knowledge Boundary

Previous works investigate the knowledge boundary to identify the known level of a knowledge piece of LLMs by quantifying the confidence or

uncertainty estimations like output consistency (Cheng et al., 2024), prompting methods (Ren et al., 2023) or knowledge probing (Ji et al., 2024). Researchers are examining the limits of parametric knowledge in LLMs with the objective of delineating the extent of the LLMs’ knowledge and iden-

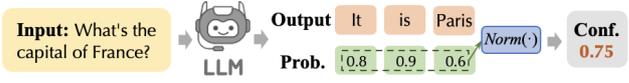
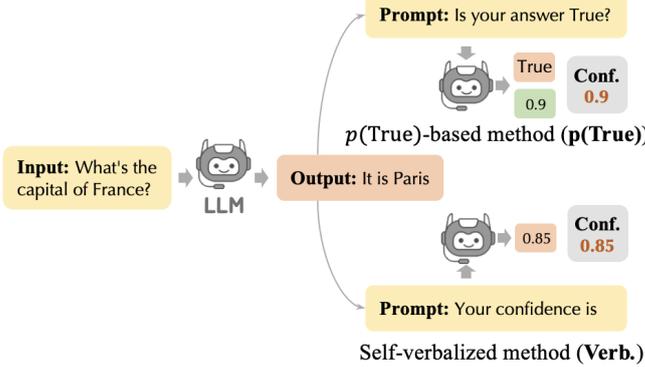
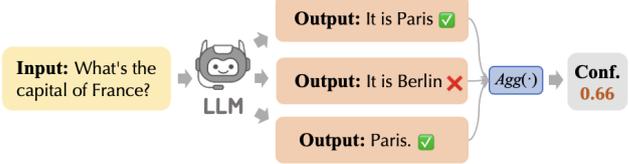
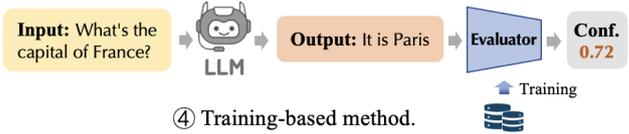
Confidence & Uncertainty Estimation Methods on LLMs	Disadvantages
 <p>① Likelihood-based method.</p>	<ul style="list-style-type: none"> <li>a. Requires normalization due to variable sequence length;</li> <li>b. Requires access to token-level probabilities, inapplicable to black-box LLMs;</li> <li>c. Fails to capture semantic meaning over token-level probabilities.</li> </ul>
 <p>② Prompting-based method.</p>	<ul style="list-style-type: none"> <li>a. Relies on prompting strategies to elicit confidence estimation, varying in different methods (Is <i>True</i> probability, numerical confidence, and word expressions, etc.);</li> <li>b. Cannot improve LLM's intrinsic confidence estimation ability.</li> <li>c. Prone to be over-confident.</li> </ul>
 <p>③ Sampling-based method.</p>	<ul style="list-style-type: none"> <li>a. Requires additional inference time cost;</li> <li>b. Varying in different aggregation methods;</li> <li>c. Cannot improve LLM's intrinsic confidence estimation ability.</li> </ul>
 <p>④ Training-based method.</p>	<ul style="list-style-type: none"> <li>a. Requires training an additional evaluator;</li> <li>b. Difficult to learn LLM's intrinsic confidence estimation on unseen domains.</li> </ul>

Figure 8: Several uncertainty estimation methods for Generative LLMs.

tifying their capability boundaries. Present studies on the knowledge boundary primarily focus on measuring the knowledge boundaries using confidence or uncertainty estimations on specialized tasks. The ambiguity of knowledge boundaries can be attributed to the knowledge distribution learned from the pre-training stage or the influence of external knowledge leading to knowledge conflict (Xu et al., 2024b) and inconsistency (Xue et al., 2023).

## G.2 Uncertainty Estimation of LLMs

To alleviate over-confidence and enhance the reliability of LLMs, reliable uncertainty estimation is essential to determine whether a question is known or not to the LLM (Geng et al., 2023). Both *Un-*

*certainty* and *Confidence* estimations can indicate the reliability degree of the responses generated by LLMs, and are generally used interchangeably (Xiao et al., 2022; Chen and Mueller, 2023; Geng et al., 2023; Lu et al., 2024). In this part, we investigate several commonly used *confidence & uncertainty* estimation methods for generative LLMs as mentioned in Sec. 5. Specifically, we denote  $\text{Conf}(x, y)$  as the confidence score associated with the output sequence  $y = [y_1, y_2, \dots, y_N]$  given the input context  $x = [x_1, x_2, \dots, x_M]$ . We also illustrate the summarized estimation methods as well as their disadvantages in Fig. 8.

**Likelihood-based Methods:** Following model calibration on classification tasks (Guo et al., 2017), Vazhentsev et al. (2023); Xue et al. (2024); Varshney et al. (2023); Wang et al. (2024) intermediately quantify sentence uncertainty over token probabilities. In traditional discriminative models, except likelihood-based methods, confidence estimations also include ensemble-based and Bayesian methods (Lakshminarayanan et al., 2017; Gal and Ghahramani, 2016; Xue et al., 2022; Wang and Yeung, 2020; Gal et al., 2016; Abdar et al., 2021), and density-based methods (Lee et al., 2018). However, this likelihood-based method requires access to token probabilities and thus being limited to white-box LLMs. The likelihood-based confidence is estimated by calculating the joint token-level probabilities over  $\mathbf{y}$  conditioned on  $\mathbf{x}$ . As longer sequences are supposed to have lower joint likelihood probabilities that shrink exponentially with length, the product of conditional token probabilities of the output should be normalized by calculating the geometric mean by the sequence length (Murray and Chiang, 2018; Malinin and Gales, 2021), and the confidence score can be represented as:

$$\text{Conf}(\mathbf{x}, \mathbf{y}) = \left( \prod_i^N p(y_i | \mathbf{y}_{<i}, \mathbf{x}) \right)^{\frac{1}{N}} \quad (9)$$

Similarly, the arithmetical average of the token probabilities is adopted in Varshney et al. (2023):

$$\text{Conf}(\mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_i^N p(y_i | \mathbf{y}_{<i}, \mathbf{x}) \quad (10)$$

Furthermore, a low probability associated with even one generated token may provide more informative evidence of uncertainty (Varshney et al., 2023). Hence, the minimum of token probabilities is also employed.

$$\text{Conf}(\mathbf{x}, \mathbf{y}) = \min \{p(y_1 | \mathbf{x}), \dots, p(y_N | \mathbf{y}_{<N}, \mathbf{x})\} \quad (11)$$

**Prompting-based Methods:** Recently, LLMs’ remarkable instruction-following ability (Brown et al., 2020) provides a view of instructing LLMs to self-estimate their confidence level to previous inputs and outputs including expressing uncertainty in words (Lin et al., 2022a; Zhou et al., 2023; Tian et al., 2023a; Xiong et al., 2024), or instructing the

LLM to self-evaluate its correctness on  $p(\text{True})$  (Kadavath et al., 2022). The  $P(\text{True})$  confidence score is implemented by simply asking the model itself if its first proposed answer  $\mathbf{y}$  to the question  $\mathbf{x}$  is true (Kadavath et al., 2022), and then obtaining the probability  $p(\text{True})$  assigned by the model, which can implicitly reflect self-reflected certainty as follows.

$$\text{Conf}(\mathbf{x}, \mathbf{y}) = p(\text{True}) = p(\mathbf{y} \text{ is True?} | \mathbf{x}) \quad (12)$$

Another method is to prompt LLMs to linguistically express tokens of confidence scores in verbalized numbers or words (Lin et al., 2022a; Mielke et al., 2022; Zhou et al., 2023; Tian et al., 2023b; Xiong et al., 2024).

The sampling-based method refers to randomly sampling multiple responses given a fixed input  $\mathbf{x}$  using beam search or temperature sampling strategies (Manakul et al., 2023; Xiong et al., 2024; Lyu et al., 2024). Various aggregation methods are adopted on sampled responses to calculate the consistency level as the confidence score. Moreover, some uncertainty quantification methods are used to calculate the entropy indicating the dispersion level of multiple outputs (Kuhn et al., 2023; Lin et al., 2023; Nikitin et al., 2024).

**Training-based Methods:** For training methods, an external evaluator trained on specific datasets is introduced to output a confidence score given an input and an output. The evaluator can be a pre-trained NLI model (Mielke et al., 2022), or a value head connected to the LLM output layer (Lin et al., 2022a; Kadavath et al., 2022), or the LLM itself (Han et al., 2024).

However, both self-verbalized and sampling methods for uncertainty estimations using extrinsic prompting or aggregation strategies with additional time costs fail to improve LLMs’ intrinsic capability of uncertainty estimation. Recent works investigate confidence learning methods to enhance the reliability of LLMs (Han et al., 2024). Li et al. (2023) introduces Inference-Time Intervention (ITI) to enhance the truthfulness of LLMs by shifting model activations during inference. Yang et al. (2023) proposes an uncertainty-aware in-context learning method leveraging uncertainty information to refine the responses but cannot improve uncertainty estimation. (Zhang et al., 2024a) proposes R-tuning to instruct LLMs to refuse unknown questions considering uncertainty estimations as binary indica-

tors. In contrast, our proposed UALIGN framework not only obtains more reliable uncertainty estimations regarding knowledge boundary information but also elicits accurate responses of LLMs.

### G.3 Factuality Alignment of LLMs

Alignment is a standard procedure to improve LLMs’ helpfulness and factuality (Bai et al., 2022a). The main goal of LLM alignment is to guide human preference through Supervised Fine-Tuning (SFT), Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022; Bai et al., 2022a) or AI feedback (Bai et al., 2022b), which may also guide LLMs to output detailed and lengthy responses (Singhal et al., 2023) but inevitably encourage hallucination. Therefore, many works explore to apply RL to improve LLMs’ factuality through Proximal Policy Optimization (PPO) (Schulman et al., 2017) with the trained reward model (Liang et al., 2024; Xu et al., 2024a) or Direct Preference Optimization (DPO) Rafailov et al. (2023) with the constructed preference dataset (Zhang et al., 2024b; Lin et al., 2024) to align with factuality preferences annotated by human beings. Xu et al. (2024a) encourage LLM to reject unknown questions using the constructed preference data by leveraging knowledge boundary feedback.

## H Experiments

### H.1 Experiments of Reliability of Uncertainty Estimations

Due to the page limitation in the main part, we present the AUROC performance results of the used confidence and entropy compared with other baseline uncertainty estimations on SciQ, NQ-Open, and LSQA as in Fig. 9, 10, and 11.

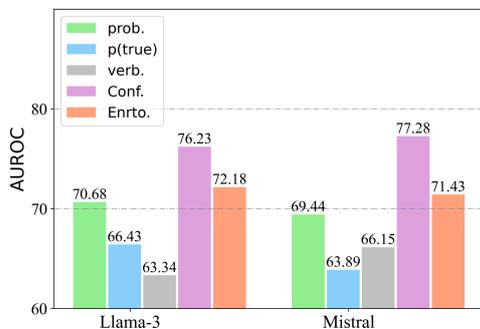


Figure 9: Results of AUORC $\uparrow$  across several confidence/uncertainty estimation methods on SciQ on Llama-3 and Mistral.

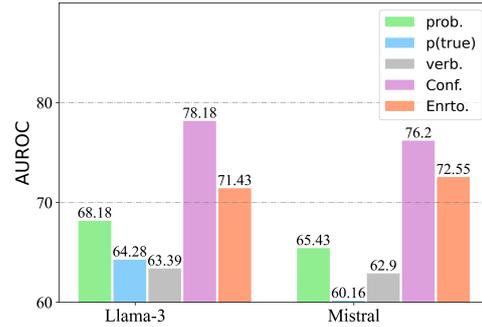


Figure 10: Results of AUORC $\uparrow$  across several confidence/uncertainty estimation methods on NQ-Open on Llama-3 and Mistral.

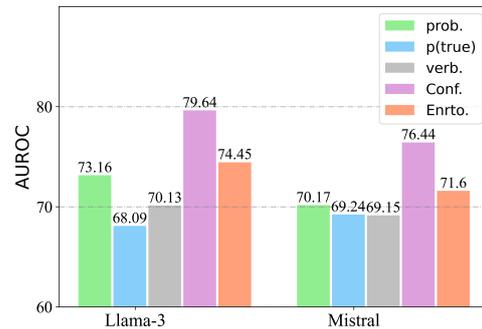


Figure 11: Results of AUORC $\uparrow$  across several confidence/uncertainty estimation methods on LSQA on Llama-3 and Mistral.

## I Few-shot Prompt Examples

10 different few-shot prompts for sampling on TriviaQA are demonstrated in Table 5.

Exemplar ID	Examples
1	Q: Which William wrote the novel Lord Of The Flies? A: Golding.
2	Q: Where in England was Dame Judi Dench born? A: York, UK.
3	Q: Neil Armstrong was a pilot in which war? A: Korean.
4	Q: How many home runs did baseball great Ty Cobb hit in the three world series in which he played? A: None.
5	Q: Who had a big 60s No 1 with Tossin' and Turnin'? A: Bobby Lewis.
6	Q: Which Disney film had the theme tune A Whole New World? A: 'Ala' ad Din.
7	Q: In basketball where do the Celtics come from? A: City of Boston.
8	Q: Which element along with polonium did the Curies discover? A: Radium.
9	Q: Who was the Egyptian king whose tomb an treasures were discovered in the Valley of the Kings in 1922? A: Tutankhamon.
10	Q: Where were the 2004 Summer Olympic Games held? A: Atina, Greece.

Table 5: Demonstrations of 1-shot examples for TriviaQA sampling to construct UALIGN dataset.

Exemplar ID	Examples
1	Q: What type of organism is commonly used in preparation of foods such as cheese and yogurt? A: mesophilic organisms.
2	Q: What phenomenon makes global winds blow northeast to southwest or the reverse in the northern hemisphere and northwest to southeast or the reverse in the southern hemisphere? A: coriolis effect.
3	Q: Changes from a less-ordered state to a more-ordered state (such as a liquid to a solid) are always what? A: exothermic.
4	Q: What is the least dangerous radioactive decay? A: alpha decay.
5	Q: Kilauea in hawaii is the world's most continuously active volcano. very active volcanoes characteristically eject red-hot rocks and lava rather than this? A: smoke and ash.
6	Q: When a meteoroid reaches earth, what is the remaining object called? A: meteorite.
7	Q: What kind of a reaction occurs when a substance reacts quickly with oxygen? A: combustion reaction.
8	Q: Organisms categorized by what species descriptor demonstrate a version of allopatric speciation and have limited regions of overlap with one another, but where they overlap they interbreed successfully? A: ring species.
9	Q: Alpha emission is a type of what? A: radioactivity.
10	Q: What is the stored food in a seed called? A: endosperm.

Table 6: Demonstrations of 1-shot examples for SciQ sampling to construct UALIGN dataset.