

# FAITH: Factuality Alignment through Integrating Trustworthiness and Honestness

Xiaoning Dong<sup>\*1,2</sup> Chengyan Wu<sup>\*3</sup> Yajie Wen<sup>3</sup> Yu Chen<sup>1</sup>  
Yun Xue<sup>2</sup> Jing Zhang<sup>†4</sup> Wei Xu<sup>‡1,2</sup> Bolei Ma<sup>5</sup>

<sup>1</sup>Tsinghua University, Institute for Interdisciplinary Information Sciences <sup>2</sup>Shanghai Qi Zhi Institute

<sup>3</sup>South China Normal University, School of Electronic Science and Engineering

<sup>4</sup>Guangzhou Richstone Data Technologies Co., Ltd. <sup>5</sup>LMU Munich & Munich Center for Machine Learning

{dongxn20, chenyu23}@mails.tsinghua.edu.cn, {chengyan.wu, yajiewen, xueyun}@m.scnu.edu.cn  
zhangjing@richstonedt.com, weixu@mail.tsinghua.edu.cn, bolei.ma@lmu.de

## Abstract

Large Language Models (LLMs) can generate factually inaccurate content even if they have corresponding knowledge, which critically undermines their reliability. Existing approaches attempt to mitigate this by incorporating uncertainty in QA prompt during training, but these numerical scores lack the semantic richness for LLM to properly understand its internal states of trustworthiness and honestness, leading to insufficient factuality alignment. We introduce **FAITH** (Factuality Alignment through Integrating Trustworthiness and Honestness), a post-training framework for factuality alignment that integrates natural-language uncertainty signals with external knowledge. Specifically, we augment training datasets by computing confidence scores and semantic entropy from LLM outputs and mapping them into a knowledge-state quadrant that describes the model’s internal knowledge possession (trustworthiness) and answering behaviors (honestness) in natural language. Based on this enhanced data, we design a reward function that considers both correctness and uncertainty signals, and fine-tune the LLM using the Proximal Policy Optimization (PPO) algorithm. To further mitigate weakly grounded responses, we design a retrieval-augmented module that retrieves relevant external passages, improving the consistency between internal and external knowledge representations. Extensive experiments on four knowledge-intensive benchmarks demonstrate that FAITH enhances the factual accuracy and truthfulness of LLMs.<sup>1</sup>

## 1 Introduction

Large Language Models (LLMs), such as GPT-4 (OpenAI, 2023), Llama3 (Dubey et al., 2024), and DeepSeek-v3 (DeepSeek-AI et al., 2024), have demonstrated impressive performance across

a broad range of natural language processing tasks. Despite these advances, growing evidence shows that LLMs may generate outputs that are fluent but factually incorrect or even fabricated, a phenomenon commonly known as hallucination (Huang et al., 2021; Ji et al., 2023). Such hallucinations pose substantial risks in knowledge-intensive and high-stakes domains, including legal, educational, and clinical applications (Alkaissi and McFarlane, 2023; Wang et al., 2023b).

A concerning type of hallucination emerges when the model possesses the necessary knowledge but fails to articulate it correctly. This disconnect between internal knowledge and external expression, often termed the *know-tell* gap (Saunders et al., 2022; Li et al., 2025), not only undermines the model’s ability to convey truthful information but also manifests as inconsistency of factual expression, where the model may produce an incorrect response in one instance yet a correct one in another (Manakul et al., 2023; Wang et al., 2023a).

In this paper, we propose to post-train LLMs for enhancing factuality. Our study identifies several limitations in recent efforts (Tian et al., 2024; Tao et al., 2024; Xue et al., 2025; Sun et al., 2025, *inter alia*): (1) while these works incorporate uncertainty into factuality alignment, they directly use the numerical values in question-answering prompts during training, which lack semantic richness and are difficult for LLMs to interpret and exploit for factuality-aligned expression; (2) they employ binary reward functions in policy learning that focus solely on whether a response is correct while overlooking the model’s internal uncertainty, or equivalently its confidence, thereby potentially encouraging guessing; and (3) they neglect the use of external knowledge, leaving potentially incorrect responses unrectified.

To address these limitations, we introduce **FAITH** (Factuality Alignment through Integrating Trustworthiness and Honestness), a post-training

<sup>\*</sup>Equal contributions.

<sup>†</sup>Corresponding authors.

<sup>1</sup>Code: <https://github.com/xndong/FAITH>

framework designed for factuality alignment in LLMs. FAITH incorporates three key designs: (1) When augmenting in-domain training datasets, beyond estimating uncertainty using consistency and semantic entropy for each question, we further map these numerical values into a knowledge-state quadrant (Liang et al., 2024) where each knowledge state is expressed in natural language and defined along two dimensions: *knowledge possession* (trustworthiness) and *answering behavior* (honestness). Unlike opaque numerical uncertainty values, incorporating knowledge states into QA prompts during training provides LLMs with semantically rich and interpretable guidance. (2) For policy optimization, we design a fine-grained reward function that considers both response correctness and model uncertainty, providing more informative feedback than a binary reward, encouraging the policy model to align its outputs with the corresponding knowledge states. (3) To further improve reliability, we construct a vector database over the Wikipedia corpus (Karpukhin et al., 2020) and train a RAG model that retrieves external knowledge from the database as contextual input to rectify potentially incorrect responses generated by the policy model.

Through FAITH, we improve LLM factuality in terms of precision and truthfulness. Extensive experiments show that FAITH consistently outperforms five recent strong baselines on three in-domain and one out-of-domain datasets. For example, on Llama3-8B, FAITH achieves 74.26% precision and 45.73% truthfulness on in-domain datasets, and 67.99% precision and 34.03% truthfulness on the out-of-domain dataset. Similar gains are observed on Mistral-7B-v0.1, demonstrating that FAITH generalizes effectively across both models and datasets.

In summary, our contributions are as follows:

1. We introduce FAITH, a novel post-training framework for factuality alignment. FAITH advances the factuality alignment by its semantically rich knowledge-state quadrant, fine-grained reward function, and employing external knowledge to ground LLM’s response.
2. We conduct extensive experiments demonstrating that FAITH consistently outperforms strong baselines, with performance gains generalizing across multiple datasets and models. Meanwhile, we provide ablations to assess the contribution of each component.

3. We provide in-depth analyses of FAITH, including the impact of different knowledge state estimation strategies on inference performance and training-time scaling behavior with varying numbers of sampled responses  $K$ .

## 2 Preliminary

| Abbr.     | Explanation                      |
|-----------|----------------------------------|
| <b>KC</b> | Known and answered correctly     |
| <b>KI</b> | Known but answered incorrectly   |
| <b>KR</b> | Known but refused to answer      |
| <b>UC</b> | Unknown but answered correctly   |
| <b>UI</b> | Unknown but answered incorrectly |
| <b>UR</b> | Unknown and refused to answer    |

Table 1: Categories of model output.

**Problem Definition.** We consider a standard *open-domain* setting, where a language model  $LLM$  is given a factual question  $q$  and generates a short-form answer  $a \sim LLM(\cdot | q)$  with probability  $LLM(a | q)$ . The answer is expected to be concise and factually correct, but in practice factuality failures may arise from uneven knowledge possession and from a gap between what the model knows and how it expresses that knowledge. However, the autoregressive generation paradigm forces the model to produce a response even when its underlying knowledge is incomplete.

Following prior work (Xue et al., 2025), we categorize model outputs into six types based on knowledge possession and response correctness, as summarized in Table 1. The research scope of this work is to encourage more instances of **KC** and **UR** (See § 4.1 for evaluation metrics).

## 3 The FAITH Framework

As shown in Figure 1, FAITH is a framework designed to improve factuality alignment for LLMs. Specifically, we first augment the training splits of three QA datasets by estimating the uncertainty of the LLM for each question and mapping the resulting numerical uncertainty values to natural-language descriptions derived from the knowledge-state quadrant (described in § 3.1). We then apply Proximal Policy Optimization (PPO) to finetune LLMs on the augmented datasets, guiding them to answer questions in a manner consistent with their knowledge states. Meanwhile, to mitigate factual failures caused by insufficient internal knowledge, we train a RAG model to rectify potentially incorrect answers produced by the policy model (de-

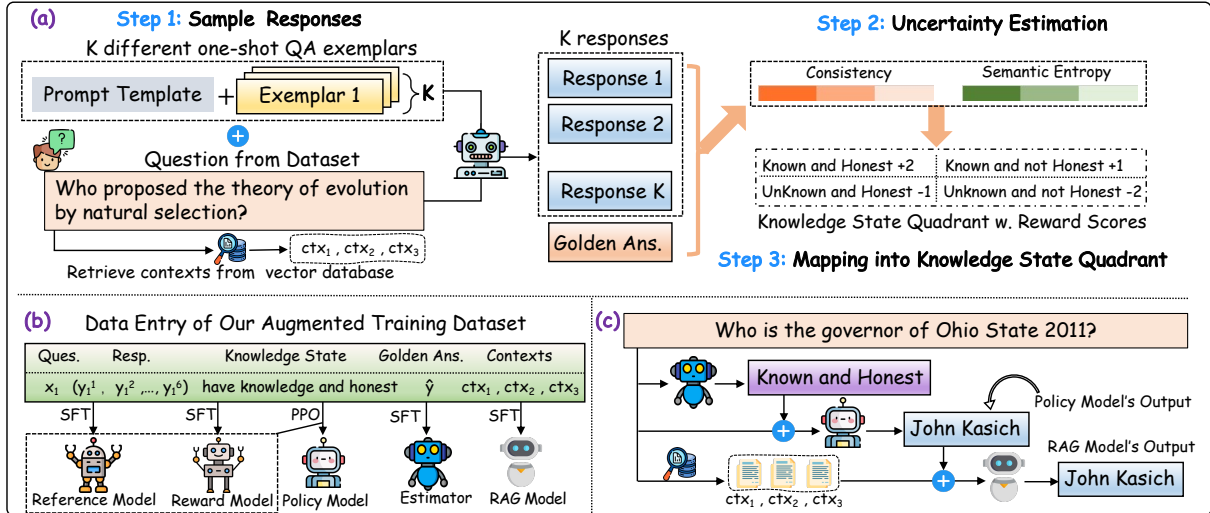


Figure 1: Illustration of the FAITH framework. Panel (a) shows the procedure for augmenting the training datasets; Panel (b) depicts model training with the augmented datasets; and Panel (c) presents the inference pipeline.

scribed in § 3.2). Finally, we describe the inference stage in § 3.3.

### 3.1 Training Dataset Augmentation

As shown in Panel (a) of Figure 1, we first sample responses for each question in the in-domain QA datasets and then estimate uncertainty, following prior work on uncertainty estimation (Xiong et al., 2024; Kuhn et al., 2023; Aichberger et al., 2025; Kang et al., 2025). We do this for two purposes: (1) to characterize the knowledge boundary of an LLM and determine whether a given question falls within it; (2) to evaluate the model’s honesty when answering that question.

#### Sampling Responses from Training Datasets.

We sample responses from the training splits of NQ-Open (Kwiatkowski et al., 2019), SciQ (Welbl et al., 2017a) and TriviaQA (Joshi et al., 2017) datasets. Each dataset  $\mathcal{D}$  contains a set of  $N$  question–answer pairs, denoted by  $\{(x_i, \hat{y}_i)\}_{i=1}^N$ , where  $x_i$  and  $\hat{y}_i$  represent the  $i$ -th question and reference answer in  $\mathcal{D}$ , respectively. Specifically, for each question  $x_i$ , we prompt the model with  $x_i$  using  $K$  different one-shot exemplars and obtain  $K$  responses, denoted by  $Y_i = \{y_i^k\}_{k=1}^K$ . We set  $K = 6$  in the main experiments, the same as that used in UAlign (Xue et al., 2025). For fair comparison, we adopt the same temperature setting ( $T=0.2$ ) as UAlign and also randomly retain half of data entries in NQ-Open and TriviaQA, while retaining all entries in SciQ.

**Uncertainty Estimation.** For uncertainty estimation in generative LLMs, we employ two measures: consistency and semantic entropy. Consistency serves as an accuracy-based proxy for confidence, reflecting the proportion of correct answers among the  $K$  generated candidate responses (Xiong et al., 2024), and is computed as follows:

$$Consistency(x_i) = \frac{1}{K} \sum_{k=1}^K \mathbb{1}\{y_i^k = \hat{y}_i\}. \quad (1)$$

Semantic entropy (SE), on the other hand, captures uncertainty from the semantic dispersion of generated answers by estimating the likelihood of each meaning cluster  $c$ , rather than each generated sequence  $y_i^k$  (Kuhn et al., 2023). It addresses the limitations of prior approaches that are often affected by response length or by semantically identical answers expressed in different surface forms. Semantic entropy is defined as:

$$SE(x_i) = - \sum_c p(c | x_i) \log p(c | x_i) \\ = - \sum_c \left( \left( \sum_{y_i^k \in c} p(y_i^k | x_i) \right) \log \left[ \sum_{y_i^k \in c} p(y_i^k | x_i) \right] \right) \quad (2)$$

At this point, we augment each data entry in  $\mathcal{D}$  from  $(x_i, \hat{y}_i)$  to  $(x_i, \hat{y}_i, Y_i, Consistency(x_i), SE(x_i))$ . However, we argue that directly using numerical uncertainty values in QA training prompts, such as “Who starred in an officer and a gentleman ### Conf: 0.833 ### Entro: -0.” in Xue et al. (2025), cannot effectively help LLMs to understand and recognize their knowledge boundaries, because raw numbers lack semantic interpretability.

**Knowledge State Mapping.** To address this limitation, we map consistency scores and semantic entropy to a knowledge-state quadrant (described below), thereby converting otherwise opaque numerical uncertainty values into semantically rich natural-language descriptions. Specifically, we characterize the knowledge states of LLMs along two dimensions: *knowledge possession* and *answering behavior*. This yields a quadrant with four knowledge states expressed in natural language: (i) **Have knowledge and honesty (KH)**, (ii) **Have knowledge but not honesty (K¬H)**, (iii) **Not have knowledge but honesty (¬KH)**, and (iv) **Not have knowledge and not honesty (¬K¬H)**.

We quantify the knowledge possession of LLMs for a given question  $x_i$  through consistency defined in Eq. 1, and the indicator function  $\mathbb{1}$  is *Positive-Recall Exact Match (PREM)*, a commonly used metric in short-form QA, where  $\text{PREM}(y_i^k, \hat{y}_i) = 1$  if  $y_i^k \in \hat{y}_i \vee \hat{y}_i \in y_i^k$ , and  $\text{PREM}(y_i^k, \hat{y}_i) = 0$  otherwise. On the other hand, we model the answering behavior of LLMs via semantic entropy.

Overall, the mapping from consistency and semantic entropy to the knowledge-state quadrant  $\mathcal{S}$  is defined as follows. We explain the details of knowledge states in Appendix A.

$$s_i = \text{KnowledgeState}(x_i) = \begin{cases} \text{KH}, & \text{if Consistency}(x_i) > 0 \text{ and } SE(x_i) = 0, \\ \text{K}\neg\text{H}, & \text{if Consistency}(x_i) > 0 \text{ and } SE(x_i) \neq 0, \\ \neg\text{KH}, & \text{if Consistency}(x_i) = 0 \text{ and } SE(x_i) = 0, \\ \neg\text{K}\neg\text{H}, & \text{otherwise.} \end{cases} \quad (3)$$

The knowledge state formulation in Eq. 3 allows us to characterize LLMs in terms of trustworthiness (knowledge possession) and honesty (answering behavior), and we finally augment each data entry from  $(x_i, \hat{y}_i)$  to  $(x_i, \hat{y}_i, Y_i, s_i)$ .

### 3.2 Training Stage of FAITH

We aim to leverage both internal and external knowledge to enhance the expression of existing knowledge, bridging the gap between knowing and telling. To this end, as shown in Panel (b) of Figure 1, we first train a policy model using PPO to align LLM responses with their internal knowledge states. We then train a RAG model to correct potentially incorrect responses by incorporating external knowledge. Finally, we introduce a knowledge state estimator that eliminates the need for sampling multiple responses during inference, thereby improving efficiency. All prompt templates used in training stage are provided in Appendix B.

**Reference Model Training.** We start from a pre-trained base model and obtain a reference model  $\pi_\mu$  through supervised fine-tuning (SFT). Specifically, the model is fine-tuned on pairs of the form  $(\text{prompt}(x_i, s_i), \hat{y}_i)$ , where the prompt incorporates both the question  $x_i$  and the knowledge state  $s_i$ , and  $\hat{y}_i$  is the reference answer. By fine-tuning the base model with these curated input-output pairs, the reference model establishes a foundation for subsequent policy optimization.

**Reward Model Training.** To align generation with knowledge states, we train a reward model with parameters  $\theta$  to evaluate generated responses together with their associated knowledge states. Unlike existing binary rewards  $r_i \in \{0, 1\}$  which focus only on whether the response is correct and ignore the model’s confidence (Yao et al., 2025; Kirichenko et al., 2025; Xue et al., 2025), we propose a fine-grained reward function that captures both response correctness and uncertainty. Specifically, we define a combined reward function:

$$R_{\text{FAITH}}(x_i, y_i^k, \hat{y}_i, s_i) = R_{\text{correctness}}(y_i^k, \hat{y}_i) + R_{\text{uncertainty}}(s_i) = \mathbb{1}_{y_i^k \equiv \hat{y}_i} + R_{\text{uncertainty}}(s_i). \quad (4)$$

where  $s_i = \text{KnowledgeState}(x_i) \in \mathcal{S}$ , and  $R_{\text{uncertainty}}(s_i) \in \{+2, +1, -1, -2\}$  is defined according to the following mapping from its knowledge state  $s_i$  in  $\mathcal{S}$ :

$$\begin{aligned} +2 &\rightarrow \text{KH}, & +1 &\rightarrow \text{K}\neg\text{H}, \\ -1 &\rightarrow \neg\text{KH}, & -2 &\rightarrow \neg\text{K}\neg\text{H}. \end{aligned}$$

We parameterize this reward function with a reward model  $RM_\theta$ . Specifically, given a dataset  $\mathcal{D}$  containing multiple tuples  $(x_i, y_i^k, \hat{y}_i, s_i, r_i^k)$ , where  $r_i^k$  is the reward value, the reward model minimizes the multi-class cross-entropy:

$$\mathcal{L}_\theta = -\mathbb{E}_{(x_i, y_i^k, \hat{y}_i, s_i, r_i^k) \sim \mathcal{D}} \left[ \log p_\theta(r_i^k | x_i, y_i^k, \hat{y}_i, s_i) \right]. \quad (5)$$

Our fine-grained reward model provides more informative feedback than binary rewards, encouraging the policy model to align its generated responses with the corresponding knowledge state, where uncertainty is expressed in natural language form rather than numeric scores.

**Policy Model Training.** Similar to reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), we employ PPO with a KL-divergence penalty to optimize LLMs for factuality alignment. Specifically, given a question  $x_i$  paired with its

knowledge state  $s_i$ , both the reference model  $\pi_\mu$  and the policy model  $\pi_\phi$  generate responses, while the reward model  $RM_\theta$  evaluates the factual reliability of a generated response  $\tilde{y}_i$  in terms of correctness and uncertainty (*i.e.*, the associated knowledge state). The training objective is to optimize  $\pi_\phi$  to maximize the expected reward:

$$\arg \max_{\pi_\phi} \mathbb{E}_{x \sim \mathcal{D}, s = \text{KnowledgeState}(x), \tilde{y} \sim \pi_\phi(x, s)} \left[ \underbrace{RM_\theta(x, \tilde{y}, s)}_{\text{reward}} - \beta \underbrace{\text{KL}[\pi_\mu(x) \parallel \pi_\phi(x, s)]}_{\text{penalty}} \right]. \quad (6)$$

**RAG Model Training.** We train a RAG model  $\pi_{rag}$  to leverage external knowledge to rectify potentially incorrect answers produced by the policy model. To this end, we first build a vector database over the Wikipedia corpus (Karpukhin et al., 2020) using the BAAI General Embedding model<sup>2</sup>. We employ IndexIVFPQ in Facebook AI Similarity Search (FAISS) (Johnson et al., 2021) as the retriever to perform similarity search. For each question  $x_i \in \mathcal{D}$ , the retriever returns the top-3 most semantically relevant passages, denoted by  $ctx_i = \{context_i^j\}_{j=1}^3$ , which are used as contextual input in the prompt. Accordingly, we augment each training entry from  $(x_i, \hat{y}_i, Y_i, s_i)$  to  $(x_i, \hat{y}_i, Y_i, s_i, ctx_i)$ . Finally, we perform retrieval-augmented fine-tuning (RAFT) (Zhang et al., 2024b) on an LLM as the rectifier, using training pairs of the form  $(\text{prompt}(x_i, s_i, \tilde{y}_i, ctx_i), \hat{y}_i)$ , where  $\tilde{y}_i$  is randomly selected from  $K$  responses in  $Y_i$ .

**Knowledge State Estimator Training.** To improve inference efficiency, we additionally train a knowledge state estimator that directly predicts the LLM’s knowledge state  $s_i$  for a given question  $x_i \in \mathcal{D}$ . Since we represent knowledge possession and answering behavior within a knowledge-state quadrant, the estimator is formulated as a four-class classification task.

Specifically, given the augmented training dataset  $\mathcal{D} = (x_i, \hat{y}_i, Y_i, s_i)_{i=1}^N$  described in § 3.1, we perform supervised fine-tuning on an LLM to serve as the knowledge-state estimator, using pairs of the form  $(\text{prompt}(x_i), s_i)$ , where the prompt incorporates the question  $x_i$  and the target label is its knowledge state  $s_i$ . The estimator is parameterized by  $\tau$ , and its training objective is defined as:

$$\mathcal{L}_\tau = -\mathbb{E}_{(x_i, s_i) \sim \mathcal{D}} [\log p_\tau(s_i | x_i)]. \quad (7)$$

<sup>2</sup>BGE-base-en-v1.5: <https://huggingface.co/BAAI/bge-base-en-v1.5>

This design enables the estimator to predict a knowledge state in a single forward pass, rather than relying on sampling  $K$  responses and computing consistency and semantic entropy. We provide empirical evaluations of the estimator’s impact on model performance in § 4.3.

### 3.3 Inference Stage of FAITH

We employ the policy model  $\pi_\phi$ , the estimator model  $Est_\tau$ , and the RAG model  $\pi_{rag}$  to perform factuality-enhanced question answering. Specifically, as shown in Panel (c) of Figure 1, given a question  $x$ , we first predict its knowledge state  $s$  in the knowledge-state quadrant using the estimator model:  $s = Est_\tau(x)$ . We then prompt the policy model  $\pi_\phi$  with  $(x, s)$  to generate the answer  $\tilde{y} = \pi_\phi(\text{prompt}(x, s))$ . Finally, we apply the RAG model to further rectify the answer produced by the policy model:  $\tilde{y}^* = \pi_{rag}(\text{prompt}(x, s, \tilde{y}, ctx))$ , obtaining the final answer  $\tilde{y}^*$ . We analyze the impact of the RAG model as a rectifier in § 4.3. All prompt templates used during inference are identical to those employed in the training stage.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets.** For training, we adopt the same widely used QA datasets as prior work to ensure a fair comparison: (1) TriviaQA (Joshi et al., 2017), where questions are from various topics and authored by trivia enthusiasts with evidence documents. (2) SciQ (Welbl et al., 2017a), which focuses on question answering in the scientific domain; and (3) NQ-Open (Kwiatkowski et al., 2019), consisting of Google search queries paired with annotated short-form answers. For evaluation, we use the test splits of these three datasets as in-domain benchmarks, and employ WebQuestions (Berant et al., 2013) as an out-of-domain dataset to assess the generalization capability of our approach. Detailed descriptions and statistics of all datasets are provided in Appendix D.

**Evaluation Metrics.** Consistent with baselines, we employ Precision (*Prec.*) and Truthfulness (*Truth.*) as evaluation metrics. Precision measures the proportion of correctly answered questions among all known questions, reflecting an LLM’s ability to accurately articulate its known knowledge. Truthfulness is defined as the proportion of correctly answered known questions plus

correctly refused unknown questions over all questions. Further details are provided in Appendix E.

**Baselines.** We compare FAITH against the following baselines spanning three categories: prompt-based, SFT-based, and RL-based methods.

(1) **ICL-CoT** (Wei et al., 2022) is a prompt-based approach that employs few-shot exemplars with chain-of-thought reasoning steps to elicit more accurate answers from LLMs without any parameter updates.

(2) **SFT** fine-tunes LLMs by minimizing the negative log-likelihood of ground-truth answers conditioned on the input questions, serving as a standard supervised learning baseline.

(3) **RL-DPO** follows Lin et al. (2024) to construct a factuality preference dataset and applies direct preference optimization to improve the factual accuracy of LLMs.

(4) **Divide-then-Align (DTA)** (Sun et al., 2025) is a multi-objective training framework for retrieval-augmented LLMs that combines DPO loss, SFT loss, and a boundary classification loss to align model behavior with explicit knowledge boundary constraints.

(5) **UAlign** (Xue et al., 2025) leverages uncertainty estimation to train LLMs to accurately express their factual knowledge boundaries, particularly for questions they cannot consistently answer correctly.

(6) **Context-DPO** (Bi et al., 2025) and (7) **In-Fact** (Cohen et al., 2025) are two additional state-of-the-art baselines included for broader comparison.

To assess the benefit of retrieval augmentation, we further implement **SFT<sub>rag</sub>** and **UAlign<sub>rag</sub>** by integrating the same RAG module used in FAITH into the respective baselines.

**Training Setup.** We implement our approach on Llama3-8B (Dubey et al., 2024) and Mistral-7B-v0.1 (Jiang et al., 2023a), using LoRA for parameter-efficient fine-tuning. Full training details are provided in Appendix C.

<sup>3</sup>For DTA with Llama3, we directly evaluate the released checkpoint, whereas for DTA with Mistral, we fine-tune Mistral-7B-Instruct on the released training data for Llama3 in a transfer setting. Meanwhile, since DTA requires augmented QA datasets with RAG context and SciQ’s augmented version was not released, its results on SciQ are unavailable (denoted as “-” in the table).

## 4.2 Main Results

We evaluate the effectiveness of our factuality alignment framework, FAITH, against strong baselines with the results shown in Table 2. The main findings are as follows:

(1) **FAITH achieves state-of-the-art performance, outperforming advanced baselines.** As shown in Table 2, FAITH consistently surpasses five baselines on three in-domain datasets and one out-of-domain dataset. For instance, on the Llama3-8B model, FAITH achieves an overall precision of 74.26% and truthfulness 45.73% on the in-domain datasets, and it attains a precision of 67.99% and a truthfulness of 34.03% on the WebQuestions dataset. We observe similar performance advantages on Mistral-7B-v0.1, except for precision on NQ-Open, demonstrating that FAITH generalizes across both models and datasets.

(2) **Natural-language knowledge states are more effective than numerical uncertainty values for guiding knowledge-boundary-aware question answering.** To assess the effectiveness of our knowledge-state-quadrant design, we compare it with numerical uncertainty values. Specifically, we construct a variant of UAlign by removing its policy optimization stage and retaining the remaining SFT stage, *i.e.*, we apply SFT with prompts that contain numerical uncertainty values. We keep all other settings unchanged. Similarly, we implement an SFT-only variant of FAITH, in which the model is prompted with natural-language knowledge states drawn from the knowledge-state quadrant. Their performance is reported in Table 2 under **UAlign<sub>sft</sub>** and **FAITH<sub>sft</sub>**, respectively.

Evaluation shows that replacing numerical uncertainty values with semantically rich knowledge states in natural language yields clear gains in guiding LLMs to recognize their knowledge boundary and answer questions accordingly. For instance, on Llama3-8B, FAITH with SFT outperforms UAlign with SFT by 2.05% in precision and 1.27% in truthfulness on average, with even larger improvements observed on Mistral-7B-v0.1. We attribute these improvements to LLMs’ preference for semantically meaningful labels (*e.g.*, “known”, “honest”) that more clearly convey knowledge boundaries. In contrast to fitting abstract numerical values, LLMs more readily interpret and leverage natural language as guidance, enabling knowledge-boundary-aware question answering.

| Method                         | TVQA (ID)        |                   | SciQ (ID)        |                   | NQ-Open (ID)     |                   | Average (ID)     |                   | WebQ-QA (OOD)    |                   |
|--------------------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|
|                                | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ |
| Llama3-8B                      |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| <b>ICL-CoT</b>                 | 66.68            | 53.37             | 72.34            | 45.90             | 57.34            | 23.60             | 65.45            | 40.95             | 65.97            | 30.85             |
| <b>SFT</b>                     | 70.80            | 52.57             | 72.18            | 45.40             | 41.41            | 16.57             | 61.46            | 38.18             | 66.46            | 31.18             |
| <b>SFT<sub>rag</sub></b>       | 71.49            | 53.04             | 73.23            | 45.97             | 42.64            | 17.58             | 62.45            | 38.86             | 66.93            | 31.85             |
| <b>RL-DPO</b>                  | 72.08            | 53.96             | 71.23            | 44.20             | 49.65            | 19.18             | 64.32            | 39.11             | 65.99            | 32.41             |
| <b>DTA<sup>3</sup></b>         | 43.99            | 31.73             | –                | –                 | 56.72            | 21.12             | –                | –                 | 61.24            | 30.71             |
| <b>UAlign</b>                  | 79.14            | 57.04             | 76.44            | 48.00             | 56.60            | 26.09             | 70.72            | 43.71             | 66.88            | 33.01             |
| <b>UAlign<sub>sft</sub></b>    | 78.76            | 56.68             | 75.87            | 47.65             | 56.02            | 25.49             | 70.22            | 43.27             | 66.12            | 32.58             |
| <b>UAlign<sub>rag</sub></b>    | 79.35            | 57.24             | 77.47            | 48.25             | 56.78            | 26.10             | 71.20            | 43.86             | 67.01            | 33.21             |
| <b>Context-DPO</b>             | 79.76            | 57.26             | 77.32            | 48.62             | 56.82            | 26.02             | 71.30            | 43.97             | 66.75            | 33.12             |
| <b>InFact</b>                  | 80.11            | 57.76             | 76.94            | 48.21             | 56.93            | 26.11             | 71.33            | 44.03             | 66.92            | 33.44             |
| <b>FAITH (ours)</b>            | <b>84.19</b>     | <b>60.69</b>      | <b>80.61</b>     | <b>49.99</b>      | <b>58.13</b>     | <b>27.58</b>      | <b>74.26</b>     | <b>45.73</b>      | <b>67.99</b>     | <b>34.03</b>      |
| <b>FAITH<sub>sft+ppo</sub></b> | 82.95            | 59.80             | 80.29            | 49.70             | 57.99            | 26.52             | 73.79            | 45.69             | 67.31            | 33.75             |
| <b>FAITH<sub>sft</sub></b>     | 81.24            | 58.77             | 78.85            | 48.65             | 56.72            | 26.21             | 72.27            | 44.54             | 66.94            | 33.10             |
| Mistral-7B-v0.1                |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| <b>ICL-CoT</b>                 | 76.73            | 54.78             | 71.87            | 44.20             | 54.47            | 18.22             | 67.69            | 39.06             | 53.43            | 35.76             |
| <b>SFT</b>                     | 74.57            | 54.77             | 65.85            | 42.50             | 50.82            | 14.42             | 63.74            | 37.08             | 52.24            | 34.33             |
| <b>SFT<sub>rag</sub></b>       | 75.87            | 55.09             | 66.08            | 43.22             | 51.03            | 14.53             | 64.33            | 37.61             | 53.14            | 34.82             |
| <b>RL-DPO</b>                  | 72.20            | 52.98             | 66.44            | 41.80             | 50.95            | 16.42             | 63.19            | 37.06             | 52.01            | 33.87             |
| <b>DTA<sup>3</sup></b>         | 41.33            | 28.78             | –                | –                 | 41.01            | 20.49             | –                | –                 | 56.53            | 23.44             |
| <b>UAlign</b>                  | 82.10            | 59.05             | 73.21            | 46.70             | 54.17            | 19.64             | 70.82            | 41.79             | 56.47            | 37.02             |
| <b>UAlign<sub>sft</sub></b>    | 81.07            | 56.47             | 72.45            | 45.87             | 43.32            | 21.55             | 65.61            | 41.30             | 55.34            | 36.87             |
| <b>UAlign<sub>rag</sub></b>    | 82.74            | 59.55             | 74.22            | 46.89             | 54.47            | 20.01             | 70.48            | 42.15             | 56.97            | 37.79             |
| <b>Context-DPO</b>             | 84.32            | 59.51             | 75.13            | 47.89             | 47.21            | 22.51             | 68.89            | 43.30             | 57.32            | 37.99             |
| <b>InFact</b>                  | 83.44            | 59.32             | 75.32            | 48.03             | 47.01            | 22.33             | 68.59            | 43.23             | 56.98            | 37.73             |
| <b>FAITH (ours)</b>            | <b>87.20</b>     | <b>60.72</b>      | 81.42            | 51.40             | <b>48.05</b>     | <b>23.91</b>      | 72.22            | <b>45.34</b>      | <b>58.04</b>     | 40.43             |
| <b>FAITH<sub>sft+ppo</sub></b> | 87.00            | 60.58             | <b>83.68</b>     | <b>51.80</b>      | 46.60            | 23.19             | <b>72.43</b>     | 45.19             | 55.48            | 38.65             |
| <b>FAITH<sub>sft</sub></b>     | 86.51            | 60.24             | 82.88            | 51.30             | 46.16            | 22.96             | 71.85            | 44.83             | 51.99            | <b>41.77</b>      |

Table 2: **Precision and Truthfulness of FAITH (ours)** vs. strong baselines on in-domain (ID) and out-of-domain (OOD) QA datasets. The Average (ID) column denotes the average performance on all three ID datasets. The subscript “sft” denotes ablation results with only supervised fine-tuning (SFT), excluding the PPO and RAG (if applicable) module. Similarly, “sft+ppo” denotes results with SFT and PPO, but excluding the RAG module. All results are reported in percentages.

(3) **Reinforcement learning with our proposed reward function improves performance.** We examine the impact of reward function design by comparing the correctness-based binary reward used in UAlign with our fine-grained reward function in Eq. 4. For example, on Llama3-8B, applying PPO with binary reward yields average gains of 0.7% in precision and 0.44% in truthfulness over SFT<sup>4</sup> on three in-domain datasets, whereas FAITH, applying PPO with our reward function, achieves larger improvements of 1.52% in precision and 1.15% in truthfulness. These results demonstrate the effectiveness of our fine-grained reward function in guiding LLM generation using both correctness and uncertainty signals.

(4) **Retrieval-Augmented Fine-Tuning aligns policy model outputs with external knowledge by rectifying potential errors.** As shown in Table 2, comparing the values under FAITH with FAITH<sub>sft+ppo</sub>, we observe consistent performance improvements across both LLMs, except for SciQ on Mistral-7B-v0.1. This demonstrates that incorporating external knowledge enhances the truthfulness of LLM responses. In addition, we manually inspect the corrections made by RAG model to the policy model outputs. Interestingly, some rectifications fail, even altering correct answers into incorrect ones, though such cases are rare. We provide an in-depth analysis of such cases in § 4.3.

### 4.3 Analysis and Discussion

**Effect of the RAG Model on Post-Hoc Correction of Policy Model Outputs.** For this analysis,

<sup>4</sup>Calculated as the difference between the metric values reported under UAlign and UAlign<sub>sft</sub> in Table 2

we conduct a statistical study on both in-domain (TriviaQA, SciQ, NQ-Open) and out-of-domain (WebQuestions) datasets, with results summarized in Figure 2. Specifically, we examine the responses produced by the policy model that are subsequently modified by the RAG model, and compute the proportion of cases where an incorrect response is corrected into a correct one versus the reverse. We find that the proportion of successful corrections consistently exceeds that of erroneous corrections across all datasets except SciQ. Notably, on TriviaQA, 87% of the policy model outputs modified by the RAG model are corrected successfully, demonstrating that effectively leveraging external knowledge can compensate for the limitations of relying solely on internal knowledge.

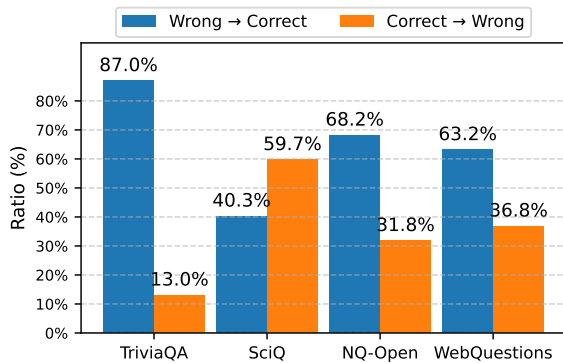


Figure 2: RAG error analysis. Each pair of bars represents the proportion of correct corrections (left) and erroneous corrections (right) among all modifications.

**Knowledge State Prediction: Estimator vs. Sample  $K$  Responses.** We further investigate the impact of different knowledge state estimation strategies on model performance during inference, comparing model-based estimation and sampling-based estimation. The model-based approach corresponds to our trained estimator model, whereas the sampling-based approach follows a two-stage pipeline: first sampling  $K$  responses using prompts with different one-shot exemplars, and then computing the knowledge state with Eq. 3. For this analysis, we set  $K = 6$ , consistent with the training stage. We present the evaluation results for precision and truthfulness in Table 3, and we observe that sampling-based estimation yields slightly higher precision and truthfulness in most cases.

These findings indicate the distribution of knowledge states can be effectively approximated by a trained LLM, while also highlighting a trade-off between efficiency and performance: sam-

pling achieves better performance and provides interpretable uncertainty measures, whereas model-based estimation avoids  $K$  rounds of inference with only minimal performance degradation.

**Training-time Scaling: The Impact of Number of Sampled Responses.** We study the training-time scaling behavior, *i.e.*, how the number of sampled responses  $K$  used during data augmentation influences the training performance of the policy model. In the main framework, we set the default value of  $K$  to 6 for efficiency. Here, specifically, we increase  $K$  from 6 to 8, 10, and 12 during data augmentation, resulting in multiple augmented datasets that differ only in the values of  $K$ . These datasets are then used to train both the estimator and the policy model. Finally, we evaluate the trained models and compare the performance. As shown in Figure 3, increasing  $K$  beyond 6 does not yield noticeable improvements in either precision or truthfulness. This suggests that sampling  $K = 6$  responses during data augmentation is already sufficient and effective to capture the distribution of the model’s knowledge states, while maintaining efficiency. In other words, while larger  $K$  values slightly expand the coverage of sampled responses, they do not translate into significant gains in downstream performance, indicating minor effects. The detailed numerical results are provided in Table 4 in the Appendix F.

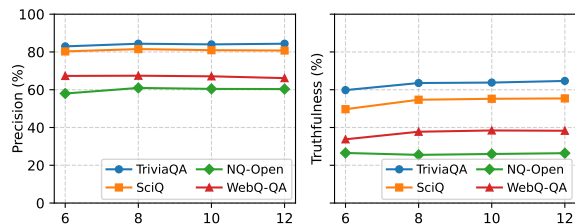


Figure 3: Training-time scaling with different numbers of sampled responses ( $K$ ) on Llama3-8B.

Finally, we present a qualitative case study in Appendix H.

## 5 Related Work

**Factuality Alignment.** To enhance LLM factuality, prior work has explored training-free strategies such as external knowledge augmentation (Kandpal et al., 2023; Jiang et al., 2023b), decoding methods (Chuang et al., 2024), and self-consistency techniques grounded in uncertainty estimation (Kadavath et al., 2022; Tian et al., 2023). More re-

| Method          | TVQA (ID)        |                   | SciQ (ID)        |                   | NQ-Open (ID)     |                   | Average (ID)     |                   | WebQ-QA (OOD)    |                   |
|-----------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|
|                 | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ |
| Llama3-8B       |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| Estimator       | 82.95            | 59.80             | 80.29            | 49.70             | 57.99            | 26.52             | 73.79            | 45.69             | 67.31            | 33.75             |
| Sampling-based  | <b>83.99</b>     | <b>59.86</b>      | <b>83.20</b>     | <b>51.50</b>      | <b>58.23</b>     | <b>26.93</b>      | <b>75.14</b>     | <b>46.10</b>      | <b>67.85</b>     | <b>34.07</b>      |
| Mistral-7B-v0.1 |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| Estimator       | 87.00            | 60.58             | 83.68            | 51.80             | <b>46.60</b>     | <b>23.19</b>      | <b>72.43</b>     | 45.19             | 55.48            | 38.65             |
| Sampling-based  | <b>87.43</b>     | <b>60.88</b>      | <b>84.01</b>     | <b>52.00</b>      | 45.66            | 22.71             | 72.37            | <b>45.20</b>      | <b>55.70</b>     | <b>38.80</b>      |

Table 3: Performance comparison between model-based and sampling-based knowledge state estimation ( $K = 6$ ). Results are reported on precision and truthfulness across all datasets and models.

cently, post-training approaches, including SFT and policy optimization, have been applied to further improve factuality (Tian et al., 2024; Tao et al., 2024; Sun et al., 2025; Xu et al., 2024b; Xue et al., 2025; Chen et al., 2025). Our work falls into this category, where we propose to map numerical uncertainty values into natural-language knowledge states and design a new reward function to enhance the model’s expression of its knowledge.

**Uncertainty Estimation.** Uncertainty Estimation (UE) has long been studied in machine learning domain, including NLP, with vast majority of previous work focusing on discriminative tasks, such as sentiment analysis (Xiao et al., 2022). Specifically, the entropy of the predictive posterior and the negative predictive posterior probability of the most probable answer are used to quantify uncertainty in predictions (Lakshminarayanan et al., 2017; Bakman et al., 2024). However, LLMs pose new challenges for uncertainty estimation due to their generative paradigm. Recent studies have extended UE to generative models, introducing heuristic or probabilistic metrics such as entropy-based scoring (Malinin and Gales, 2021), semantic entropy that accounts for meaning equivalence in multiple generations (Kuhn et al., 2023), and similarity-based methods applied in tasks like machine translation (Fomicheva et al., 2020; Lin et al., 2022b). Other works explore black-box UE by leveraging sampled outputs (Chen and Mueller, 2024; Manakul et al., 2023), or prompt-based approaches where models verbalize their own confidence (Kadavath et al., 2022). Training-based methods have also been proposed to enhance linguistic self-assessments of uncertainty (Lin et al., 2022a).

**Retrieval-Augmented Generation.** Retrieval-Augmented Generation (RAG) has become a widely adopted strategy for mitigating the limita-

tions of parametric knowledge in LLMs by grounding responses in external evidence (Guu et al., 2020; Izacard and Grave, 2021; Zhong et al., 2023). Unlike knowledge editing methods that directly modify model parameters (Yao et al., 2023; Xu et al., 2024a), RAG enables models to dynamically access updated information without retraining, reducing the risk of catastrophic forgetting. However, recent studies reveal that contradictions may arise when retrieved knowledge conflicts with internal representations (Mei et al., 2024; Ni et al., 2024). To mitigate this, retrieval-based methods increasingly leverage structured repositories such as knowledge graphs for more reliable grounding (Baek et al., 2023; Zhang et al., 2024a). Other approaches enhance factuality by incorporating retrieval into inference-time interventions, including memory-augmented architectures (Li et al., 2022) and entity-level embedding integration (Kang et al., 2022; de Jong et al., 2022). Generally, RAG provides a scalable way to continuously integrate knowledge, offering advantages over task-specific parameter editing.

## 6 Conclusion

We present a post-training framework, called FAITH, for factuality alignment in LLMs. Our approach estimates uncertainty and translates the numerical values into natural-language knowledge states that measure the knowledge possession and answering behavior of LLMs. Meanwhile, we design a fine-grained reward function to incentivize both correctness and uncertainty of LLM’s response. Finally, we introduce a trained RAG model to rectify potentially incorrect responses generated by policy model. Experiments show that FAITH substantially outperforms recent baselines in truthfulness and precision. We hope this work contributes to building more faithful and factual LLMs as part of the broader community effort.

## Limitations

**Reward Function Design.** Our reward function is derived from heuristic rules that are straightforward to formulate and intuitively easy to interpret. In practice, we observe that this design works well empirically and provides meaningful guidance for aligning model behavior. However, the current formulation lacks rigorous theoretical guarantees, leaving room for future work to establish a stronger theoretical foundation for its effectiveness.

**Computational Overhead.** During dataset construction, we sample  $K$  responses and build a vector database. At inference time, our pipeline first uses  $Est_\tau$  to estimate the knowledge state  $s$ , then applies the policy model  $\pi_\phi$  to generate an answer, and finally employs  $\pi_{rag}$  for rectification. Even without rectification, two model inferences are required, rather than a single end-to-end pass. For future work, we plan to explore more efficient approaches for cognitive-state estimation, such as lightweight estimators derived from LLM internal representations (Zhu et al., 2025).

**Unexplored Aspects of RAG Effectiveness.** Our current study does not investigate how the quality of the data used to build the vector database affects FAITH’s performance. For example, on the SciQ dataset with Mistral-7B, incorporating external knowledge does not improve correction effectiveness, which may be related to the quality of the retrieved context. In addition, we have not explored more effective ways of leveraging external knowledge for rectification, such as integrating RAFT directly the SFT stage and accordingly applying PPO training on top of the RAFT-enhanced model, rather than additionally train a RAG model.

## Ethical Considerations

This research was conducted in accordance with the ACM Code of Ethics. The datasets (Joshi et al., 2017; Welbl et al., 2017a; Kwiatkowski et al., 2019; Berant et al., 2013) that we use are publicly available. We report only aggregated results in the main paper. We did not share any Personally Identifiable Data with this paper. The project itself does not raise any risks.

## Acknowledgments

CW and YX are supported in part by the Guangdong Basic and Applied Basic Research Foundation under Grant 2023A1515011370, the

National Natural Science Foundation of China (32371114), the Characteristic Innovation Projects of Guangdong Colleges and Universities (No. 2018KTSCX049), and the Guangdong Provincial Key Laboratory (No. 2023B1212060076). XD and WX are supported in part by the National Key R&D Program of China 2023YFC3304802 and National Natural Science Foundation of China (NSFC) Grant U2268202 and 62176135.

## References

- Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. 2025. [Improving uncertainty estimation through semantically diverse language generation](#). In *The Thirteenth International Conference on Learning Representations*.
- Hussam Alkaissi and Samy I McFarlane. 2023. Artificial hallucinations in chatgpt: implications in scientific writing. *Cureus*, 15(2).
- Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. [Knowledge-augmented language model prompting for zero-shot knowledge graph question answering](#). In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 78–106, Toronto, Canada. Association for Computational Linguistics.
- Yavuz Faruk Bakman, Duygu Nur Yaldiz, Baturalp Buyukates, Chenyang Tao, Dimitrios Dimitriadis, and Salman Avestimehr. 2024. [MARS: Meaning-aware response scoring for uncertainty estimation in generative LLMs](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7752–7767, Bangkok, Thailand. Association for Computational Linguistics.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on Freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.
- Baolong Bi, Shaohan Huang, Yiwei Wang, Tianchi Yang, Zihan Zhang, Haizhen Huang, Lingrui Mei, Junfeng Fang, Zehao Li, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, and Shenghua Liu. 2025. [Context-DPO: Aligning language models for context-faithfulness](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10280–10300, Vienna, Austria. Association for Computational Linguistics.
- Jiuhai Chen and Jonas Mueller. 2024. [Quantifying uncertainty in answers from any language model and enhancing their trustworthiness](#). In *Proceedings of*

- the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), *ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 5186–5200. Association for Computational Linguistics.
- Xilun Chen, Iliia Kulikov, Vincent-Pierre Berges, Barlas Oğuz, Rulin Shao, Gargi Ghosh, Jason Weston, and Wen tau Yih. 2025. [Learning to reason for factuality](#). *Preprint*, arXiv:2508.05618.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2024. [Dola: Decoding by contrasting layers improves factuality in large language models](#). In *The Twelfth International Conference on Learning Representations*.
- Roi Cohen, Russa Biswas, and Gerard de Melo. 2025. [InFact: Informativeness alignment for improved LLM factuality](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 17876–17888, Suzhou, China. Association for Computational Linguistics.
- Michiel de Jong, Yury Zemlyanskiy, Nicholas FitzGerald, Fei Sha, and William W. Cohen. 2022. [Mention memory: incorporating textual knowledge into transformers through entity mention attention](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, and Chenggang et al. Zhao. 2024. [DeepSeek-V3 Technical Report](#). *Preprint*, arXiv:2412.19437.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, and Artem et al. Korenev. 2024. [The Llama 3 Herd of Models](#). *Preprint*, arXiv:2407.21783.
- Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. [Unsupervised quality estimation for neural machine translation](#). *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. [Retrieval augmented language model pre-training](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3929–3938. PMLR.
- Yi-Chong Huang, Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2021. [The factual inconsistency problem in abstractive text summarization: A survey](#). *CoRR*, abs/2104.14839.
- Gautier Izacard and Edouard Grave. 2021. [Leveraging passage retrieval with generative models for open domain question answering](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pages 874–880. Association for Computational Linguistics.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023a. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023b. [Active retrieval augmented generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 7969–7992. Association for Computational Linguistics.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. [Billion-scale similarity search with gpus](#). *IEEE Trans. Big Data*, 7(3):535–547.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, and 17 others. 2022. [Language models \(mostly\) know what they know](#). *CoRR*, abs/2207.05221.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. [Large language models struggle to learn long-tail knowledge](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 15696–15707. PMLR.
- Minki Kang, Jinheon Baek, and Sung Ju Hwang. 2022. [KALA: knowledge-augmented language model adaptation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for*

- Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 5144–5167. Association for Computational Linguistics.
- Zhewei Kang, Xuandong Zhao, and Dawn Song. 2025. [Scalable best-of-n selection for large language models via self-certainty](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Polina Kirichenko, Mark Ibrahim, Kamalika Chaudhuri, and Samuel Bell. 2025. [Abstentionbench: Reasoning LLMs fail on unanswerable questions](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. [Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation](#). In *The Eleventh International Conference on Learning Representations*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, 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. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. [Simple and scalable predictive uncertainty estimation using deep ensembles](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Siheng Li, Cheng Yang, Taiqiang Wu, Chufan Shi, Yuji Zhang, Xinyu Zhu, Zesen Cheng, Deng Cai, Mo Yu, Lemao Liu, Jie Zhou, Yujie Yang, Ngai Wong, Xixin Wu, and Wai Lam. 2025. [A survey on the honesty of large language models](#). *Transactions on Machine Learning Research*. Survey Certification.
- Zonglin Li, Ruiqi Guo, and Sanjiv Kumar. 2022. [Decoupled context processing for context augmented language modeling](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 21698–21710. Curran Associates, Inc.
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024. [Learning to trust your feelings: Leveraging self-awareness in LLMs for hallucination mitigation](#). In *Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP*, pages 44–58, Bangkok, Thailand. Association for Computational Linguistics.
- Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Scott Yih, and Xilun Chen. 2024. [FLAME : Factuality-aware alignment for large language models](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a. [Teaching models to express their uncertainty in words](#). *Transactions on Machine Learning Research*.
- Zi Lin, Jeremiah Zhe Liu, and Jingbo Shang. 2022b. [Towards collaborative neural-symbolic graph semantic parsing via uncertainty](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 4160–4173, Dublin, Ireland. Association for Computational Linguistics.
- Andrey Malinin and Mark Gales. 2021. [Uncertainty estimation in autoregressive structured prediction](#). In *International Conference on Learning Representations*.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. [Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 9004–9017. Association for Computational Linguistics.
- Lingrui Mei, Shenghua Liu, Yiwei Wang, Baolong Bi, and Xueqi Cheng. 2024. [SLANG: new concept comprehension of large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 12558–12575. Association for Computational Linguistics.
- Shiyu Ni, Keping Bi, Jiafeng Guo, and Xueqi Cheng. 2024. [When do llms need retrieval augmentation? mitigating llms’ overconfidence helps retrieval augmentation](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 11375–11388. Association for Computational Linguistics.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural*

- Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.*
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. 2022. [Self-critiquing models for assisting human evaluators](#). *CoRR*, abs/2206.05802.
- Xin Sun, Jianan Xie, Zhongqi Chen, Qiang Liu, Shu Wu, Yuehe Chen, Bowen Song, Zilei Wang, Weiqiang Wang, and Liang Wang. 2025. [Divide-then-align: Honest alignment based on the knowledge boundary of RAG](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11461–11480, Vienna, Austria. Association for Computational Linguistics.
- Shuchang Tao, Liuyi Yao, Hanxing Ding, Yuexiang Xie, Qi Cao, Fei Sun, Jinyang Gao, Huawei Shen, and Bolin Ding. 2024. [When to trust llms: Aligning confidence with response quality](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 5984–5996. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2024. [Fine-tuning language models for factuality](#). In *The Twelfth International Conference on Learning Representations*.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 5433–5442. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023a. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. [Self-instruct: Aligning language models with self-generated instructions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13484–13508. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017a. [Crowdsourcing multiple choice science questions](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017b. [Crowdsourcing multiple choice science questions](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text, NUT@EMNLP 2017, Copenhagen, Denmark, September 7, 2017*, pages 94–106. Association for Computational Linguistics.
- Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2022. [Uncertainty quantification with pre-trained language models: A large-scale empirical analysis](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 7273–7284. Association for Computational Linguistics.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie Fu, Junxian He, and Bryan Hooi. 2024. [Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs](#). In *The Twelfth International Conference on Learning Representations*.
- Derong Xu, Ziheng Zhang, Zhihong Zhu, Zhenxi Lin, Qidong Liu, Xian Wu, Tong Xu, Wanyu Wang, Yuyang Ye, Xiangyu Zhao, Enhong Chen, and Yefeng Zheng. 2024a. [Editing factual knowledge and explanatory ability of medical large language models](#). In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM 2024, Boise, ID, USA, October 21-25, 2024*, pages 2660–2670. ACM.
- Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai Fan, Lu Chen, and Kai Yu. 2024b. [Rejection improves reliability: Training LLMs to refuse unknown questions using RL from knowledge feedback](#). In *First Conference on Language Modeling*.
- Boyang Xue, Fei Mi, Qi Zhu, Hongru Wang, Rui Wang, Sheng Wang, Erxin Yu, Xuming Hu, and Kam-Fai Wong. 2025. [UAlign: Leveraging uncertainty estimations for factuality alignment on large language models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6002–6024, Vienna, Austria. Association for Computational Linguistics.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. [Editing large language models: Problems, methods, and opportunities](#). In *Proceedings of the 2023 Conference on Empirical Methods in*

*Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 10222–10240. Association for Computational Linguistics.

Zijun Yao, Yantao Liu, Yanxu Chen, Jianhui Chen, Junfeng Fang, Lei Hou, Juanzi Li, and Tat-Seng Chua. 2025. *Are reasoning models more prone to hallucination?* Preprint, arXiv:2505.23646.

Shuo Zhang, Liangming Pan, Junzhou Zhao, and William Yang Wang. 2024a. *The knowledge alignment problem: Bridging human and external knowledge for large language models*. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2025–2038, Bangkok, Thailand. Association for Computational Linguistics.

Tianjun Zhang, Shishir G Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E. Gonzalez. 2024b. *RAFT: Adapting language model to domain specific RAG*. In *First Conference on Language Modeling*.

Zexuan Zhong, Zhengxuan Wu, Christopher D. Manning, Christopher Potts, and Danqi Chen. 2023. *Mquake: Assessing knowledge editing in language models via multi-hop questions*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 15686–15702. Association for Computational Linguistics.

Yubo Zhu, Dongrui Liu, Zecheng Lin, Wei Tong, Sheng Zhong, and Jing Shao. 2025. *The LLM already knows: Estimating LLM-perceived question difficulty via hidden representations*. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 1160–1176, Suzhou, China. Association for Computational Linguistics.

## A Derivation of Rules for Mapping Uncertainty Values to Knowledge States

The proposed rule-based mapping from uncertainty values (*Consistency* and *SE*) to natural-language knowledge states is as follow:

Let  $Y_i = \{y_i^k\}_{k=1}^K$  be  $K$  sampled responses from model and  $\hat{y}$  be the reference answer. *Knowledge possession* is captured by *Consistency*, where:

- $Consistency > 0$  indicates there exists at least one response  $y_i^k$  matches the ground-truth answer  $\hat{y}$ , meaning the model possesses the required knowledge with a very high probability.
- $Consistency = 0$  represents that the model fails to answer the question correctly with  $K$  times, indicating the model does not possess relevant knowledge to the question. For a

given question, we assume that if the LLM possesses the relevant knowledge, the probability of answering it correctly is  $p_\theta = 0.5$ . We set the confidence level to  $\alpha = 0.05$ . If none of the  $K$  sampled responses are correct, then with confidence  $1 - \alpha$  we can reject the hypothesis that the model’s probability of producing a correct answer satisfies  $p_\theta \geq 0.5$ . Under the assumption  $p_\theta = 0.5$ , sampling  $K = 6$  responses is sufficient to show that if none of them are correct, the model is unlikely to possess the relevant knowledge.

*Answering Behavior* is measured by *semantic entropy*:

- The magnitude of semantic entropy reflects the model’s uncertainty at the semantic level: a higher value indicates diverse or conflicting semantic outputs (greater ambiguity), while a lower value suggests more consistent and deterministic semantic interpretations.
- Semantic entropy equals zero when all generated outputs are semantically equivalent, *i.e.*, they fall into the same semantic cluster with no competing interpretations. From the semantic perspective, the model is completely certain, exhibiting neither ambiguity nor polysemy.

Given these interpretations, the mapping rules follow a logically consistent decision path:

1. If  $Consistency > 0$  and  $SE = 0$ , the model is judged to possess the relevant knowledge of a question and honestly provides consistent correct responses, corresponding to the knowledge state **KH**.
2. If  $Consistency > 0$  and  $SE \neq 0$ , the model produces a mix of correct and incorrect answers, indicating insufficient mastery of the knowledge to express it accurately. The reason for this gap could be decoding strategy, hallucination snowballing, misalignment issues (Liang et al., 2024). This corresponds to the knowledge state **K¬H**.
3. If  $Consistency = 0$  and  $SE = 0$ , the model lacks correct knowledge but converges on a single interpretation, corresponding to the knowledge state **¬KH**.
4. In all other cases, the knowledge state is classified as **¬K¬H**.

Overall, the mapping is determined by two factors:

$$\underbrace{\text{Knowledge possession (Consistency)}}_{\text{know}}$$

and

$$\underbrace{\text{Answer honesty (Semantic Entropy)}}_{\text{tell}}$$

which together define a quadrant of four cognitive states, ensuring both interpretability and completeness.

## B Prompt Templates

We illustrate the prompt templates used in this work in Figure 4, detailing the input structure, incorporated knowledge states, and output format. The templates explicitly define how a question is combined with its corresponding knowledge state, optionally with retrieved external context, and then formatted to elicit model responses. By making this structure explicit, the figure clarifies how prompts guide the model during both training and inference, ensuring consistency across stages. Moreover, the design rationale highlights how natural-language descriptions of knowledge states are integrated into the prompt, which is essential for conveying uncertainty information in a semantically interpretable way.

## C Details of Training

All experiments are conducted on a cluster equipped with  $4 \times$  NVIDIA A40 GPUs.

For supervised fine-tuning (SFT) of both the reference model and estimator in FAITH, we train for 3 epochs. We adopt the Adam optimizer with an initial learning rate of  $2e-4$ . We apply LoRA with a rank of 32, alpha of 16, and a dropout rate of 0.05, targeting all layers. The batch size per device is set to 8, with gradient accumulation steps of 8, leading to a total batch size of 256. The learning rate scheduler follows a cosine decay with a warmup ratio of 0.0.

For policy optimization, both the reward model (RM) and the PPO stages are trained for 2 epochs. We adopt the Adam optimizer with an initial learning rate of  $1e-5$ . LoRA is applied with a rank of 8, alpha of 16, and a dropout rate of 0.05, again targeting all layers. The per-device batch size is set to 4, with gradient accumulation steps of 8, leading to a total batch size of 128. The learning rate

### Prompt template for sampling responses

```
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 ###: {input_question}
### Answer ###:
```

### Prompt template for supervised fine-tuning of the reference model

```
You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.
### Question ###: {THE QUESTION FROM DATASET}
### Self-Eval ###: {THE KNOWLEDGE STATE FROM DATASET}
### Output ###: {REFERENCE ANSWER}
```

### Prompt template for policy model optimization

```
You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.
### Question ###: {THE QUESTION FROM DATASET}
### Self-Eval ###: {THE KNOWLEDGE STATE FROM DATASET}
### Answer ###: {REFERENCE ANSWER}
```

### Prompt template for RAFT a RAG model

```
You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.
### Question ###: {THE QUESTION FROM DATASET}
### Self-Eval ###: {THE KNOWLEDGE STATE FROM DATASET}
### Prior Judgment ###: {RANDOMLY SELECTED RESPONSE FROM  $Y_i$ }
### Retrieve Documents ###: related passages: ###passage 1###,###passage 2###,###passage 3###
### Posterior Answer ###: {REFERENCE ANSWER}
```

### Prompt template for RAG model in inference

```
You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.
### Question ###: {THE QUESTION FROM DATASET}
### Self-Eval ###: {THE KNOWLEDGE STATE FROM DATASET}
### Prior Judgment ###: {POLICY MODEL'S OUTPUT}
### Retrieve Documents ###: related passages: ###passage 1###,###passage 2###,###passage 3###
### Posterior Answer ###: {REFERENCE ANSWER}
```

### Prompt template for supervised fine-tuning of the estimator model.

```
You are an excellent Question-Answering assistant. Please answer the following question based on your knowledge.
### Question ###: {THE QUESTION FROM DATASET}
### Self-Eval ###: {THE KNOWLEDGE STATE FROM DATASET}
```

Figure 4: All the prompt templates employed in FAITH.

scheduler is cosine decay, and the warmup ratio is 0.0, consistent with the SFT stage.

## D Details of Dataset

**SciQ:** The SciQ dataset (Welbl et al., 2017b) contains 13,679 crowdsourced science examination questions covering subjects such as physics, chemistry, and biology. Although originally released in multiple-choice format, in our setting all answer options are removed, and each question is reformulated as an open-ended query requiring a direct answer. For most questions, an accompanying paragraph with supporting evidence is provided, offering factual context that can be utilized to guide answer generation and factual alignment. In our experiments, 11,679 samples are used for training and 1,000 samples are reserved for validation, with the remaining questions serving as an in-domain test set.

**TriviaQA:** TriviaQA (Joshi et al., 2017) is a large-scale reading comprehension dataset containing over 650K question-answer-evidence triples, with questions authored by trivia enthusiasts and evidence documents collected from Wikipedia and the web. In our work, to construct the augmented dataset, we pre-process and sample half of the original training set.

**NQ-Open:** NQ-Open (Kwiatkowski et al., 2019) is an open-domain QA benchmark derived from the Natural Questions dataset, where real user queries are paired with English Wikipedia passages as the knowledge source. In our work, we employ NQ-Open for augmented dataset construction. Similarly, to ensure fair comparison and reduce computational cost, we sample half of the original training data.

**Web-Questions:** The WebQuestions dataset (Berant et al., 2013) comprises 6,642 question-answer pairs, where each question can be answered using Freebase, a large-scale knowledge graph. The majority of questions are centered around a single named entity and reflect typical queries collected from the web around 2013. In our experiments, we only employ its test set (1348 item) for the evaluation under the out-of-domain evaluation setting.

## E Details of Evaluation Metrics

Truthfulness quantifies the proportion of correct responses among all provided answers, reflecting the

LLM’s overall reliability in expressing knowledge. The formula for Truthfulness is given as follows:

$$\text{Truthfulness} = \frac{\text{UR} + \text{KC}}{\text{KC} + \text{KI} + \text{KR} + \text{UC} + \text{UI} + \text{UR}} \quad (8)$$

Precision measures the proportion of correctly answered questions among those for which the model possesses the relevant knowledge, reflecting the LLM’s ability to accurately convey known facts. The formula for Precision is given as follows:

$$\text{Precision} = \frac{\text{KC}}{\text{KI} + \text{KC} + \text{KR}} \quad (9)$$

## F Numerical Results of Training-time Scaling

Table 4 reports the detailed numerical results corresponding to the training-time scaling analysis. The table compares precision and truthfulness across different values of  $K$  (6, 8, 10, 12). Consistent with Figure 3, the results show that increasing  $K$  beyond 6 does not yield noticeable gains, confirming that  $K = 6$  is sufficient to capture the model’s knowledge state distribution during data augmentation.

## G Contribution of the Estimator

In Table 2, we report the results of FAITH under “FAITH (ours)”. To further examine the contribution of the estimator within the full FAITH framework, we additionally compare the performance of FAITH with model-based and sampling-based knowledge-state estimation, respectively, under the setting  $K = 6$ . As shown in Table 5, although the sampling-based method yields a slightly higher value in both precision and truthfulness, our estimator model achieves comparable performance (within approximately 1%) while reducing inference latency by avoiding multiple sampling rounds. These observations are consistent with the findings in Table 3.

## H Case Study

In our method FAITH, one focus is to train RAG model to align the policy model’s output with external knowledge. The RAG model is provided with retrieved passages as context, allowing it to rectify or retain the policy model’s responses. In this section, we analyze three types of corrections, with representative cases shown in Table 6, 7, and 8, as case studies. Specifically, the three correction types are summarized as follows:

| # Responses | TVQA (ID)        |                   | SciQ (ID)        |                   | NQ-Open (ID)     |                   | Average (ID)     |                   | WebQ-QA (OOD)    |                   |
|-------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|
|             | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ |
| Llama-3-8B  |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| <b>K=6</b>  | 82.95            | 59.80             | 80.29            | 49.70             | 57.99            | 26.52             | 73.79            | 45.69             | 67.31            | 33.75             |
| <b>K=8</b>  | 84.36            | 63.55             | 81.52            | 54.70             | 60.95            | 25.51             | 75.61            | 47.92             | 67.42            | 37.76             |
| <b>K=10</b> | 84.01            | 63.79             | 80.94            | 55.20             | 60.44            | 25.98             | 75.13            | 48.32             | 67.10            | 38.43             |
| <b>K=12</b> | 84.36            | 64.65             | 80.76            | 55.40             | 60.35            | 26.40             | 75.16            | 48.82             | 66.15            | 38.28             |

Table 4: Training-time scaling with different numbers of sampled responses ( $K$ ) on Llama-3-8B.

| Method                                | TVQA (ID)        |                   | SciQ (ID)        |                   | NQ-Open (ID)     |                   | Average (ID)     |                   | WebQ-QA (OOD)    |                   |
|---------------------------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|------------------|-------------------|
|                                       | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ | Prec. $\uparrow$ | Truth. $\uparrow$ |
| Llama3-8B                             |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| <b>FAITH<sub>sampling-based</sub></b> | <b>85.02</b>     | <b>60.98</b>      | <b>81.43</b>     | <b>50.33</b>      | <b>59.25</b>     | <b>28.02</b>      | <b>75.23</b>     | <b>46.44</b>      | <b>68.67</b>     | <b>34.89</b>      |
| <b>FAITH<sub>Estimator</sub></b>      | 84.19            | 60.69             | 80.61            | 49.99             | 58.13            | 27.58             | 74.26            | 45.73             | 67.99            | 34.03             |
| Mistral-7B-v0.1                       |                  |                   |                  |                   |                  |                   |                  |                   |                  |                   |
| <b>FAITH<sub>sampling-based</sub></b> | <b>88.12</b>     | <b>61.43</b>      | <b>82.53</b>     | <b>52.23</b>      | <b>48.97</b>     | <b>24.59</b>      | <b>73.21</b>     | <b>46.08</b>      | <b>58.87</b>     | <b>41.23</b>      |
| <b>FAITH<sub>Estimator</sub></b>      | 87.20            | 60.72             | 81.42            | 51.40             | 48.05            | 23.91             | 72.22            | 45.34             | 58.04            | 40.43             |

Table 5: Performance comparison of FAITH<sub>sampling-based</sub> and FAITH<sub>Estimator</sub> on in-domain and out-of-domain datasets.

- 1. Implicitly Supported Correction:** The initial answer from the policy model was incorrect, but after applying our trained RAG model, the final answer was corrected. Notably, the retrieved passages did not verbatim reproduce the correct answer, but contained key information or semantic cues related to the correct answer. Details can be found in Table 6.
- 2. Explicitly Supported Correction:** The policy model initially produced an incorrect output, but after applying our trained RAG model, the final output was corrected. In this process, the retrieved content from RAG not only directly reproduced the correct answer but also provided additional information related to it, thereby supporting the model’s correction. Details can be found in Table 7.
- 3. Misleading Override:** The policy model initially produced the correct answer. However, after applying our trained RAG model, the output was incorrectly altered. This occurred because the retrieved content contained misleading information that contradicted the correct answer, ultimately leading to an erroneous output. Details can be found in Table 8.

## I Reproducibility Statement

We have made several efforts to ensure the reproducibility of our work. We provided detailed descriptions of the datasets used in our experiments, all of which are publicly available. Our

method is thoroughly explained in dedicated sections § 3, and we also provide detailed training parameters § 4.1. Finally, we have publicized the code at <https://github.com/xndong/FAITH>.

We hope that these measures will facilitate the replication of our work by other researchers and further advance the field.

## J The Use of Large Language Model

The authors acknowledge the use of OpenAI ChatGPT solely for enhancing the coherence of the final manuscript, and providing assistance with coding for data processing.

| Type                            | Question & Answer   | Retrieved Passages  | Our Analysis  |
|---------------------------------|---|---|---|
| Implicitly Supported Correction | <p><b>Q:</b> Protists play critically important ecological roles as producers and, on the other end of food webs, as what?</p> <p><b>A:</b> decomposers<br/> <b>A1:</b> consumers<br/> <b>A2:</b> decomposers</p> | <p><b>I.</b> In real world ecosystems, there is more than one food chain for most organisms, since most organisms eat more than one kind of food or are eaten by more than one type of predator. A diagram that sets out the intricate network of intersecting and overlapping food chains for an ecosystem is called its food web. Decomposers are often left off food webs, but if included, they mark the end of a food chain. <b>Thus, food chains start with primary producers and end with decay and decomposers.</b></p> <p><b>II.</b> Food webs have trophic levels and positions. Basal species, such as plants, form the first level and are the resource-limited species that feed on no other living creature in the web. Basal species can be autotrophs or detritivores, including decomposing organic material and its associated microorganisms, which we defined as detritus, micro-inorganic material, and associated microorganisms (MIP), and vascular plant material.</p> <p><b>III.</b> The microbial food web refers to the combined trophic interactions among microbes in aquatic environments. These microbes include viruses, bacteria, algae, and heterotrophic protists (such as ciliates and flagellates). In aquatic environments, microbes constitute the base of the food web. Single-celled photosynthetic organisms such as diatoms and cyanobacteria are generally the most important primary producers in the open ocean. Many of these cells, especially cyanobacteria, are too small to be captured and consumed by small crustaceans and planktonic larvae. Instead, these cells are consumed by phagotrophic protists, which are readily consumed by larger organisms.</p> | <p>Some protists do function as A1 (consumers), but the model fails to accurately address the specific context of “the other end of food webs” posed in the question. This suggests the model’s insufficiency in effectively utilizing its internal knowledge to answer the question, or a failure to correctly understand the question, particularly its key constraints. After introducing the retrieved information to our trained RAG model, it successfully corrected the answer to A2 (decomposers). The retrieved information, in paragraph I, provides the crucial background knowledge: “Thus food chains start with primary producers and end with decay and decomposers.” This information does not explicitly state that “protists are decomposers.” Instead, it requires the model to synthesize this information with “protists” and “the other end of food webs” to deduce the correct answer. We denote this process as <b>Implicitly Supported Correction</b>.</p> |

Table 6: Case Study Analysis 1 of Model Responses. Specifically, **A** denotes Ground Truth, **A1** denotes the policy model’s output, and **A2** denotes the final output.

| Type                            | Question & Answer   | Retrieved Passages   | Our Analysis  |
|---------------------------------|---|--|---|
| Explicitly Supported Correction | <p><b>Q:</b> Rita Coolidge sang the title song for which Bond film?</p> <p><b>A:</b> Octopussy</p> <p><b>A1:</b> North by Northwest</p> <p><b>A2:</b> Octopussy</p> | <p><b>I.</b> Octopussy is the soundtrack for the eponymous thirteenth James Bond film. The score was composed by John Barry, the lyrics by Tim Rice. <b>The opening theme, All Time High is sung by Rita Coolidge and is one of six Bond film title songs or songs that are not named after the film’s title.</b> The original compact disc released in 1985 by A&amp;M Records, was recalled because of a printing error and became a rarity.</p> <p><b>II.</b> Another Way to Die is a song by American musicians Jack White and Alicia Keys. Written and produced by White as the theme song to the 2008 James Bond film Quantum of Solace, it was released as a single in the United States on September 30, 2008, and in Europe on October 20, 2008.</p> <p><b>III.</b> Tomorrow Never Dies is the song, performed by Sheryl Crow, which served as the theme song to the James Bond film of the same name. The song was co-written by Crow and the song’s producer Mitchell Froom, and became her fifth UK Top 20 hit, peaking at No. 12 in 1997. Another song, Tomorrow Never Dies, written by the movie’s composer David Arnold and performed by k.d. lang, was originally produced as the official theme tune. When Crow’s song became the official theme, the k.d. lang song was relegated to the end credits, and renamed Surrender.</p> | <p>The model’s initial response, A1 (North by Northwest), is a significant factual error, as this film is not even part of the James Bond series. This indicates a substantial knowledge gap or a “hallucination” in the model’s internal knowledge base. After the RAG intervention, the model successfully corrected the answer to A2 (Octopussy). The retrieved information in paragraph I contains all the key details required to rectify the error. The passage explicitly states, “The opening theme, All Time High is sung by Rita Coolidge” and that the soundtrack was for the film Octopussy. The model simply needed to match the key entity from the question, “Rita Coolidge”, with the retrieved text to directly find the name of the film for which she sang the theme. The entire process involves direct information extraction and localization, requiring almost no complex reasoning. We denote this process as <b>Explicitly Supported Correction</b>.</p> |

Table 7: Case Study Analysis 2 of Model Responses. Specifically, **A** denotes Ground Truth, **A1** denotes the policy model’s output, and **A2** denotes the final output.

| Type                | Question & Answer   | Retrieved Passages   | Our Analysis   |
|---------------------|---|--|--|
| Misleading Override | <p><b>Q:</b> Which grand slam did Pete Sampras not win in the 20th century?<br/> <b>A:</b> French<br/> <b>A1:</b> French Open<br/> <b>A2:</b> Wimbledon</p> | <p><b>I.</b> As the Swiss national anthem played Federer was overcome with emotion after finally capturing the elusive title at Roland Garros. This match was momentous in the history of tennis. After missing the chance to equal Pete Sampras' then-record of fourteen Grand Slam championships of all time when he lost to Rafael Nadal in the final of the Australian Open earlier in the year, Federer finally did so by winning the French Open for the first time. Sampras himself commented on Federer following the victory saying, Regardless he [Federer] goes down as the greatest ever.</p> <p><b>II.</b> In the eight Wimbledons inclusive between 1993 and 2000, 1996 was the only year that Sampras would fail to win the championship at Wimbledon. Sampras lost in the quarterfinals of Wimbledon to the eventual winner, Richard Krajicek, the tournament's 17th-seed. The match lasted three long sets, with Krajicek winning 7-5, 7-6, 6-4. In the quarterfinals of the US Open, Sampras vomited on the court at 1-1 in the final set tiebreak (due to dehydration) while facing Àlex Corretja; nonetheless, Sampras would win that match.</p> <p><b>III.</b> He beat former champion Michael Stich in the fourth round and met Sampras in the quarterfinals. By that time, he had managed to turn his notably weak slice backhand into an aggressive top-spin shot. Krajicek shocked the tennis world by defeating Sampras in straight sets, becoming the only player to beat Sampras in a Wimbledon singles match in the eight-year period from 1993 until Sampras' fourth-round loss to Roger Federer in the 2001 tournament.</p> | <p>This is a failure case of a 'correct-to-incorrect' reversal caused by the RAG. The model's initial judgment, A1 (French Open), was correct, indicating that its internal knowledge base already contained the key fact about Sampras's career. However, the intervention of RAG instead led to a degradation in performance. The core of the failure lies in the Retrieval stage. The retrieved information, though related to the key entities "Pete Sampras" and "Grand Slam", did not align with the question's specific requirement ("did not win" in the 20th century). The retrieved content, particularly in paragraphs II and III, repeatedly and in detail described a specific loss Sampras had at Wimbledon (in 1996 to Krajicek). Phrases like "fail to win the championship at Wimbledon" became a strong and irrelevant distracting signal. When generating the final answer, the model over-relied on this incorrectly retrieved and distracting content, thereby ignoring its own correct prior knowledge. It was misled into outputting the incorrect answer A2 (Wimbledon). We denote this process as <b>Misleading Override</b>.</p> |

Table 8: Case Study Analysis 3 of Model Responses. Specifically, **A** denotes Ground Truth, **A1** denotes the policy model's output, and **A2** denotes the final output.