

KnowTuning: Knowledge-aware Fine-tuning for Large Language Models

Anonymous ACL submission

Abstract

Despite their success at many natural language processing (NLP) tasks, large language models (LLMs) still struggle to effectively leverage knowledge for knowledge-intensive tasks, manifesting limitations such as generating incomplete, non-factual, or illogical answers. These limitations stem from inadequate knowledge awareness of LLMs during vanilla fine-tuning. To address these problems, we propose a knowledge-aware fine-tuning (KnowTuning) method to improve fine-grained and coarse-grained knowledge awareness of LLMs. We devise a fine-grained knowledge augmentation stage to train LLMs to identify difficult fine-grained knowledge in answers. We also propose a coarse-grained knowledge comparison stage to train LLMs to distinguish between reliable and unreliable knowledge, in three aspects: completeness, factuality, and logicity. Extensive experiments on both generic and medical question answering (QA) datasets confirm the effectiveness of KnowTuning, through automatic and human evaluations, across various sizes of LLMs. We further verify that KnowTuning generates more facts with less factual error rate under fine-grained facts evaluation.

1 Introduction

Large language models (LLMs) have become a default solution for many natural language processing (NLP) scenarios, including the question answering (QA) task (Brown et al., 2020; Ouyang et al., 2022; Qin et al., 2023). To achieve strong performance, most LLM first accumulate substantial knowledge by pre-training on extensive datasets (Jiang et al., 2023; Touvron et al., 2023). Then, in the supervised fine-tuning (SFT) stage, these LLMs further learn downstream domain knowledge and how to exploit the corresponding knowledge to answer diverse questions (Wei et al., 2022; Chung et al., 2022; Wang et al., 2023f; Peng et al., 2023; Kang et al., 2023; Wang et al., 2023c).

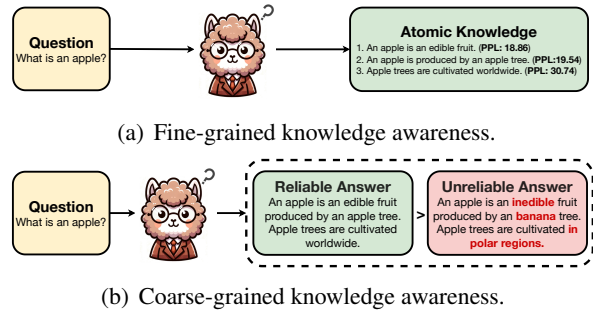


Figure 1: Illustrations of vanilla fine-tuned LLMs lacking knowledge awareness. (a) Vanilla fine-tuned LLMs struggles to identify the fine-grained knowledge to answer a specific question precisely. (b) Vanilla fine-tuned LLMs cannot effectively distinguish between reliable knowledge and unreliable knowledge in answers.

However, fine-tuned LLMs often struggle to effectively leverage knowledge for complex knowledge-intensive question-answering (Yu et al., 2023a; Bai et al., 2023; Chen et al., 2023b; Chang et al., 2023). Concretely, many recent studies indicate that LLMs are susceptible to generating incomplete answers, offering incomprehensive and insufficient knowledge (Singhal et al., 2022; Bian et al., 2024; Xu et al., 2023a); non-factual answers, delivering factually incorrect knowledge (Wang et al., 2023a; Min et al., 2023; Wang et al., 2023b); or illogical answers, providing incoherent and poorly structured knowledge (Chen et al., 2023b; Zhong et al., 2023; Kang et al., 2023). Although recent method FactTune (Tian et al., 2023) improves the factuality of answers by increasing the proportion of correct facts, it ignores other critical aspects, such as completeness (Min et al., 2023) and logicity (Xu et al., 2023a).

We hypothesize that these limitations of LLMs arise from insufficient fine-grained and coarse-grained knowledge awareness during vanilla fine-tuning (Bian et al., 2024; Ji et al., 2023; Dou et al., 2023; Hua et al., 2024). On the one hand, as illustrated in Figure 1, at the fine-grained level, vanilla fine-tuned LLMs face difficulties in identifying de-

tailed atomic knowledge within the answer, leading to inadequate awareness of fine-grained knowledge. On the other hand, at the coarse-grained level, LLMs frequently fail to distinguish between reliable and unreliable knowledge in answers, indicating a lack of coarse-grained knowledge awareness. Consequently, there is a pressing need for designing knowledge-aware fine-tuning methods. This leads to our central research question: *how can we effectively improve both the fine-grained and coarse-grained knowledge awareness of LLMs to address complex knowledge-intensive tasks?*

To this end, we propose a novel knowledge-aware fine-tuning method, named KnowTuning, which aims to improve the fine-grained and coarse-grained knowledge awareness of LLMs. KnowTuning consists of two stages: (i) fine-grained knowledge augmentation, and (ii) coarse-grained knowledge comparison. In the first stage, we filter difficult atomic knowledge with high perplexity from original answers, and rewrite fine-grained QA pairs based on the filtered knowledge. After that, we subsequently use both the original and fine-gained QA pairs to train LLMs. In the second stage, we adopt several knowledge-disturbing techniques to construct coarse-grained knowledge comparison sets along three dimensions, completeness, factuality, and logicity. Specifically, we generate answers that are worse in terms of completeness, factuality, or logicity, by deleting, revising, and shuffling the atomic knowledge. Besides, we rephrase original answers based on the atomic knowledge to prevent overfitting. Finally, we combine the rephrased answers and answers with worse completeness, factuality, and logicity as our knowledge comparison sets. We adopt direct preference optimization (DPO) (Rafailov et al., 2023) for optimizing LLMs on our coarse-grained knowledge comparison sets.

We conduct experiments on a generic QA dataset and a medical QA dataset using automatic and human evaluations. Experimental results demonstrate the effectiveness of our proposed method KnowTuning, assessing completeness, factuality, and logicity across various sizes of LLMs. Furthermore, we demonstrate that KnowTuning not only generates more facts but also reduces the factual error rate during fine-grained facts evaluation.

In summary, our main contributions are:

- We focus on improving the fine-grained and coarse-grained knowledge awareness of LLMs via fine-tuning for knowledge-intensive tasks.

- We introduce KnowTuning, a novel method that fine-tunes LLMs to leverage fine-grained knowledge augmentation and coarse-grained knowledge comparison to improve fine-grained and coarse-grained knowledge awareness of LLMs.
- We demonstrate the effectiveness of KnowTuning in the generic and medical domain QA datasets through automatic and human evaluations, across various sizes of LLMs. Furthermore, KnowTuning generates more facts with less factual error rate under fine-grained facts evaluation.

2 Related work

2.1 LLMs for knowledge-intensive Tasks

Large language models (LLMs) have been applied to various knowledge-intensive tasks (Moi-seev et al., 2022; Yu et al., 2023b; Khattab et al., 2022; Tian et al., 2023; Zhang et al., 2023a; Xu et al., 2023b; Mishra et al., 2023; Nguyen et al., 2023). Previous work mainly focus on knowledge-intensive tasks with short-form answers. Liu et al. (2022b) use few-shot demonstrations to elicit relevant knowledge statements from LLMs for QA tasks. Liu et al. (2022a) train a neural model to generate relevant knowledge through reinforcement learning for QA tasks. Liu et al. (2023a) propose a unified model for generating relevant knowledge and solving QA tasks.

However, these methods primarily address multiple-choice QA, rather than the more complex open-ended knowledge-intensive QA tasks (Krishna et al., 2021; Kadavath et al., 2022; Liu et al., 2022a, 2023a; Kang et al., 2023), which aim to solve questions that require detailed explanations and extensive domain knowledge. Recent research indicates that LLMs face challenges in tackling complex knowledge-intensive QA tasks (Yu et al., 2023a; Bai et al., 2023; Chang et al., 2023). In particular, they are prone to generating responses that are non-factual (Lee et al., 2022; Sun et al., 2023; Su et al., 2022), incomplete (Singhal et al., 2022; Bian et al., 2024), or illogical (Chen et al., 2023b; Zhong et al., 2023). Recently, for open-ended knowledge-intensive tasks, Tian et al. (2023) propose a method FacTune to improve factuality. Specifically, they first automatically evaluate the proportion of correct facts in candidate answers as factuality scores, and fine-tuning LLMs to increase the likelihood of generating answers with higher factuality scores. In contrast, our work focus on

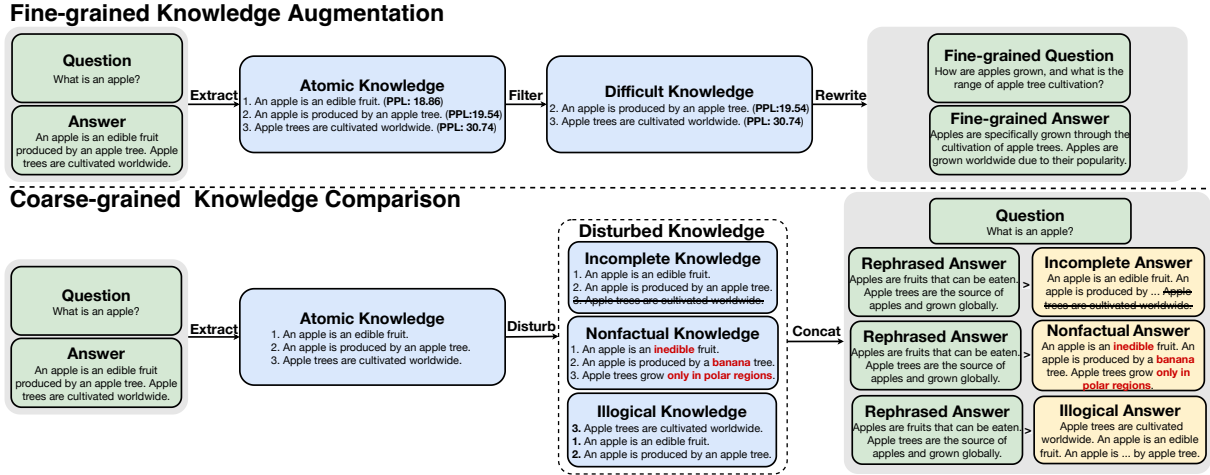


Figure 2: Overview of KnowTuning. KnowTuning leverages fine-grained knowledge augmentation and coarse-grained knowledge comparison to improve the knowledge awareness of LLMs.

improving the knowledge awareness of LLMs at multiple essential aspects simultaneously, for solving complex knowledge-intensive QA tasks.

2.2 Fine-tuning for LLMs

Fine-tuning is a kind of method to optimize pre-trained LLMs for further learning downstream domain knowledge and how to exploit the corresponding knowledge to answer diverse questions (Brown et al., 2020; Ouyang et al., 2022). Previously, fine-tuning is mainly focused on enhancing general-purpose QA abilities of LLMs (Wang et al., 2022; Wei et al., 2022; Longpre et al., 2023). These approaches mainly adopt human-annotated datasets to build the QA dataset. Recently, an alternative strategy involves generating QA datasets through the utilization of advanced LLMs to create answers to a variety of questions (Wang et al., 2023f; Shumailov et al., 2023).

Another line of fine-tuning methods fuse information about the quality of the generated answers into the supervision signals (Zhao et al., 2023; Guo et al., 2023; Wang et al., 2023d; Dong et al., 2023; Chen et al., 2024). Rafailov et al. (2023) propose direct preference optimization (DPO) to directly optimize LLMs on the pair-wise comparison set. Song et al. (2023) propose Preference Ranking Optimization (PRO) to fine-tune LLMs on list-wise comparison sets. Yuan et al. (2023) propose a margin-rank loss to optimize the LLMs on comparison sets. Since collecting large-scale human judgment for the quality of generated answers is expensive, Bai et al. (2022) and Lee et al. (2023) propose reinforcement learning from AI feedback (RLAIF) methods to leverage off-the-shelf LLMs to annotate

general helpfulness scores. In contrast, our work focuses on enhancing the fine-grained and coarse-grained knowledge-awareness of LLMs to improve performance in terms of completeness, factuality, and logicity simultaneously.

3 Method

In this section, we detail the KnowTuning method. First, we introduce the preliminaries. Then, we introduce the fine-grained knowledge augmentation. Next, we introduce coarse-grained knowledge comparison in detail. Finally, a training process for KnowTuning is explained.

3.1 Preliminaries

Supervised fine-tuning. Supervised fine-tuning (SFT) aims to train pre-trained LLMs to understand and answer natural language questions. Formally, given a QA dataset $\mathcal{D} = \{(q_i, a_i)\}_{i=1}^N$, where q_i and a_i denotes a question and a corresponding answer. The training objective of SFT is to minimize the following loss:

$$\mathcal{L}_{\text{SFT}} = - \sum_{j=1}^{|a_i|} \log P_{\pi_{\text{SFT}}}(a_{i,j} | a_{i,<j}, q_i), \quad (1)$$

where $a_{i,j}$ denotes the j -th token of a_i .

Atomic Knowledge. Since individual facts can well cover the knowledge in answers (Nenkova and Passonneau, 2004; Zhang and Bansal, 2021; Liu et al., 2023b; Min et al., 2023; Wei et al., 2024), we break an answer into individual facts as atomic knowledge. The atomic knowledge is a short statement conveying one piece of fact, which is a more

233 fine-grained unit than a sentence. Specifically, we
 234 extract atomic knowledge set \mathcal{K} from the original
 235 answers a as follows:

$$236 \quad \mathcal{K}_i = \{k_i^j\}_{j=1}^{|\mathcal{K}_i|} = \text{Extract}(a_i), \quad (2)$$

237 where $\text{Extract}(\cdot)$ is implemented by prompting
 238 OpenAI models to extract atomic knowledge, fol-
 239 lowing Min et al. (2023).

240 3.2 Fine-grained Knowledge Augmentation

241 As illustrated in Figure 2, to improve the fine-
 242 grained knowledge awareness of LLMs, we filter
 243 difficult atomic knowledge for LLMs, and rewrite
 244 fine-grained QA pairs based on the difficult knowl-
 245 edge. After that, we subsequently use both the
 246 original and fine-gained QA pairs to train LLMs.
 247 To filter the difficult atomic knowledge for LLMs,
 248 we first compute the generation perplexity ppl_i^j of
 249 each atomic knowledge k_i^j conditioned on q_i as
 250 follows:

$$251 \quad ppl_i^j = \sqrt[n]{\frac{1}{\sum_{m=1}^{|k_i^j|} P_{\pi_{SFT}}(k_{i,m}^j | k_{i,<m}^j, q_i)}}}. \quad (3)$$

252 Since high perplexity ppl indicates the lack of
 253 knowledge awareness of LLMs on specific atomic
 254 knowledge, we select α percent of the atomic
 255 knowledge set \mathcal{K}_i in descending order of perplexity
 256 to form the difficult knowledge set \mathcal{K}_i^* . Then, we
 257 rewrite the question q_i as a fine-grained question
 258 q_i^* relevant to difficult knowledge \mathcal{K}_i^* , as follows:

$$259 \quad q_i^* = \text{Rewrite}(q_i, \mathcal{K}_i^*), \quad (4)$$

260 where $\text{Rewrite}(\cdot)$ is implemented by prompting
 261 OpenAI models. In addition, we rewrite the answer
 262 based on the difficult knowledge set as the fine-
 263 grained answer:

$$264 \quad a_i^* = \text{Rewrite}(\mathcal{K}_i^*). \quad (5)$$

265 Finally, we combine the original QA dataset \mathcal{D}
 266 and the fine-grained QA pairs as the fine-grained
 267 knowledge augmentation dataset \mathcal{D}_{ka} as:

$$268 \quad \mathcal{D}_{ka} = \mathcal{D} \cup \{q_i^*, a_i^*\}_{i=1}^N. \quad (6)$$

269 3.3 Coarse-grained Knowledge Comparison

270 To improve coarse-grained knowledge awareness
 271 of LLMs in terms of completeness, factuality and
 272 logicity, we construct three comparison sets by
 273 deleting, revising, and shuffling atomic knowledge.

Knowledge completeness comparison. To im-
 274 prove knowledge completeness awareness of
 275 LLMs, we construct the knowledge completeness
 276 comparison set by randomly deleting the atomic
 277 knowledge. Specifically, we first randomly delete
 278 atomic knowledge k in the atomic knowledge set
 279 \mathcal{K} as incomplete knowledge set:
 280

$$281 \quad \mathcal{K}_i^c = \text{Delete}(\mathcal{K}_i), \quad (7)$$

282 where $\text{Delete}(\cdot)$ refers to randomly delete β per-
 283 cent of atomic knowledge k . Then, we concate-
 284 nate leftover atomic knowledge of the incomplete
 285 knowledge set as an incomplete answer:

$$286 \quad a_i^c = \text{Concat}(\mathcal{K}_i^c). \quad (8)$$

287 In addition, to avoid overfitting on the original an-
 288 swers (Jain et al., 2023), we rephrase the original
 289 answers based on the original atomic knowledge
 290 set as:

$$291 \quad a_i^r = \text{Rewrite}(\mathcal{K}_i). \quad (9)$$

292 Finally, we combine the rephrased answer a_i^r and
 293 the incomplete answer a_i^c into knowledge complete-
 294 ness comparison set as follows:

$$295 \quad \mathcal{D}_{kcc} = \{(q_i, (a_i^r, a_i^c))\}_{i=1}^N. \quad (10)$$

Knowledge factuality comparison. To improve
 296 the knowledge factuality awareness of LLMs, we
 297 construct the knowledge factuality comparison set
 298 by revising the atomic knowledge as nonfactual
 299 atomic knowledge. Specifically, we first revise the
 300 atomic knowledge set \mathcal{K}_i as follows:
 301

$$302 \quad \mathcal{K}_i^f = \text{Revise}(\mathcal{K}_i), \quad (11)$$

303 where $\text{Revise}(\cdot)$ is implemented by prompting Ope-
 304 nAI models to revise the atomic knowledge to the
 305 wrong atomic knowledge. Then, we concatenate
 306 all atomic knowledge in the nonfactual knowledge
 307 set as:

$$308 \quad a_i^f = \text{Concat}(\mathcal{K}_i^f). \quad (12)$$

309 Finally, we combine the rephrased answer a_i^r and
 310 the nonfactual answer a_i^f into knowledge factuality
 311 comparison set as follows:

$$312 \quad \mathcal{D}_{kfc} = \{(q_i, (a_i^r, a_i^f))\}_{i=1}^N. \quad (13)$$

Knowledge logicity comparison. To improve
 313 the knowledge logicity awareness of LLMs, we
 314 construct the knowledge logicity comparison
 315 set by randomly shuffling the atomic knowledge.
 316

Specifically, we first randomly shuffle all atomic knowledge in the atomic knowledge set \mathcal{K} as the illogical knowledge set:

$$\mathcal{K}_i^l = \text{Shuffle}(\mathcal{K}_i), \quad (14)$$

where $\text{Shuffle}(\cdot)$ is implemented by shuffling the order of all atomic knowledge k in the atomic knowledge set \mathcal{K} . Then, we follow the shuffled order to concatenate all atomic knowledge in the illogical knowledge set as an illogical answer:

$$a_i^l = \text{Concat}(\mathcal{K}_i^l). \quad (15)$$

Finally, we combine the rephrased answer a_i^r and the illogical answer a_i^l into knowledge logicity comparison set as follows:

$$\mathcal{D}_{klc} = \{(q_i, (a_i^r, a_i^l))\}_{i=1}^N. \quad (16)$$

Finally, we combine the knowledge completeness comparison set, the knowledge factuality comparison set, and the knowledge logicity comparison set as the coarse-grained knowledge comparison set:

$$\mathcal{D}_{kc} = \mathcal{D}_{kcc} \cup \mathcal{D}_{kfc} \cup \mathcal{D}_{klc}. \quad (17)$$

3.4 Training

To improve the knowledge awareness of LLMs for solving complex knowledge-intensive tasks, KnowTuning includes fine-grained knowledge augmentation training and coarse-grained knowledge comparison training. Specifically, we first train LLMs on fine-grained knowledge augmentation dataset \mathcal{D}_{ka} , resulting in a model denoted as π_{ka} . Then, KnowTuning aims to further improve the coarse-grained knowledge awareness of the model π_{ka} in completeness, factuality, and logicity. To accomplish this, we rewrite the DPO (Rafailov et al., 2023) loss to obtain the coarse-grained knowledge comparison loss as follows:

$$\mathcal{L}_{kc} = \mathbb{E}_{(q, (a_w, a_l)) \sim \mathcal{D}_{kc}} \left[\log \sigma \left(\beta \log \frac{\pi_{kc}(a_w|q)}{\pi_{ka}(a_w|q)} \right) - \beta \log \frac{\pi_{kc}(a_l|q)}{\pi_{ka}(a_l|q)} \right], \quad (18)$$

where (a_w, a_l) denotes the answer pair of the question $q \in \mathcal{D}_{kc}$, and a_w is the better answer.

4 Experiments

4.1 Research questions

We aim to answer the following research questions in our experiments: **RQ1**: How does KnowTuning

perform on generic and medical QA under automatic evaluation and human evaluation? **RQ2**: How does KnowTuning perform on generic and medical QA under fine-grained facts evaluation? **RQ3**: How do fine-grained knowledge augmentation and coarse-grained knowledge comparison affect the performance of KnowTuning?

4.2 Datasets

We conduct experiments on general domain and domain-specific knowledge-intensive question-answering datasets:

- **Dolly** (Conover et al., 2023) is a general domain QA dataset carefully curated by thousands of human annotators. Since we focus on open-ended generic domain QA, we filter QA pairs of “open_qa” and “general_qa” categories.
- **MedQuAD** (Abacha and Demner-Fushman, 2019) is a medical domain QA dataset, which is collected from 12 National Institutes of Health websites. Following August et al. (2022), we filter QA pairs of the category “Information” for giving detailed information about medical terms. More details of datasets are in Appendix A.

4.3 Baselines

We compare our model with the following baselines:

- **Base** denotes that testing Llama2-base models (Touvron et al., 2023) under zero-shot setting.
- **SFT** (Ouyang et al., 2022) represents vanilla fine-tuning backbone LLMs on QA datasets according to Eq. 1.
- **RLAIF** (Bai et al., 2022; Lee et al., 2023) leverages LLMs to annotate overall helpfulness scores for candidate answers, and construct overall helpfulness comparison sets based on the scores.
- **FactTune** (Tian et al., 2023) constructs factuality comparison sets by calculating the proportion of correct facts in candidate answers.

More details of baselines are in Appendix B.

4.4 Evaluation Metrics

We present our experimental results using two evaluation metrics: automatic evaluation and human-based evaluation. Since ROUGE (ROUGE, 2004) and BLEU (Papineni et al., 2002) can not effectively evaluate the quality of answers for complex questions (Krishna et al., 2021; Xu et al., 2023a), recent studies propose to use GPT-4 for evaluating the quality of LLMs answers (Zheng et al., 2024; Dubois et al., 2023; Fu et al., 2023). Consequently,

Method	Dataset	Completeness			Factuality			Logicity			Avg. gap
		Win	Tie	Lose	Win	Tie	Lose	Win	Tie	Lose	
Backbone Language Model: Llama2-7b-base											
KnowTuning vs Base	Dolly	88.50*	3.00	8.50	73.00*	20.00	7.00	80.50*	12.00	7.50	+73.00
KnowTuning vs SFT		78.50*	5.50	16.00	37.00*	46.50	16.50	50.50*	34.00	15.50	+39.33
KnowTuning vs RLAIIF		69.50*	5.00	25.50	32.00*	49.00	19.00	46.50*	39.00	14.50	+29.67
KnowTuning vs FactTune		64.50*	10.00	25.50	30.00*	53.00	17.00	31.50*	56.00	13.00	+23.50
KnowTuning vs Base	MedQuAD	93.00*	3.00	4.00	72.50*	20.50	7.00	85.00*	8.50	6.50	+77.67
KnowTuning vs SFT		81.00*	3.50	15.50	46.50*	37.50	16.00	64.50*	21.50	14.00	+48.83
KnowTuning vs RLAIIF		85.00*	2.50	12.50	41.00*	38.50	20.50	50.50*	30.00	19.50	+41.33
KnowTuning vs FactTune		83.00*	3.50	13.50	40.50*	36.50	23.00	50.50*	31.50	18.00	+39.83
Backbone Language Model: Llama2-13b-base											
KnowTuning vs Base	Dolly	85.50*	6.50	8.00	66.00*	24.50	9.50	81.00*	13.00	6.00	+69.67
KnowTuning vs SFT		77.00*	5.00	18.00	35.50*	49.50	15.00	45.00*	40.00	15.00	+36.50
KnowTuning vs RLAIIF		73.50*	4.00	22.50	33.50*	52.50	14.00	46.50*	40.50	13.00	+34.67
KnowTuning vs FactTune		68.50*	6.50	25.00	30.50*	55.00	14.50	36.00*	54.00	10.00	+28.50
KnowTuning vs Base	MedQuAD	92.50*	2.50	5.00	73.50*	17.50	9.00	84.00*	8.00	8.00	+76.00
KnowTuning vs SFT		86.50*	3.50	10.00	45.50*	41.00	13.50	60.00*	31.00	9.00	+53.16
KnowTuning vs RLAIIF		82.50*	5.00	12.50	38.50*	48.00	13.50	54.00*	38.50	7.50	+47.17
KnowTuning vs FactTune		78.00*	4.50	17.50	37.00*	47.00	16.00	48.50*	39.50	12.00	+39.33

Table 1: Main results on generic QA and medical QA datasets evaluated by GPT-4. The scores marked with * mean KnowTuning outperforms the baseline significantly with p -value < 0.05 (sign. test), following Guan et al. (2021).

given golden label as a reference, we employ GPT-4 to rate generated answers on three aspects: completeness, factuality, and logicity, on a range of 1 to 10. Following Singhal et al. (2022); Zheng et al. (2024); Zhang et al. (2023b), we define completeness, factuality and logicity as: (i) **Completeness**: it examines whether the answers provide comprehensive and sufficient knowledge to the questions. (ii) **Factuality**: it examines whether the knowledge in the answers is factually correct. (iii) **Logicity**: it examines whether the knowledge in the answers is logically structured. Following Li et al. (2023); Chen et al. (2023a), we define “Win-Tie-Lose” as: (i) **Win**: KnowTuning wins twice, or wins once and ties once. (ii) **Tie**: KnowTuning ties twice, or wins once and loses once. (iii) **Lose**: KnowTuning loses twice, or loses once and ties once.

We also employ human judgments as the gold standard for assessing the quality of answers. Specifically, human evaluators perform pair-wise comparisons of the top-performing models identified in automatic evaluations. They are presented with a question with a golden answer, and asked to judge two generated answers on three aspects: completeness, factuality, and logicity.

To evaluate the capabilities of LLMs at a fine-grained level, we follow Min et al. (2023) to conduct fine-grained facts evaluation. Specifically, we first break candidate answers into individual facts, and use *gpt-3.5-turbo* to measure the correctness of each fact based on the golden answer as a reference.

Following Tian et al. (2023), we report the number of correct facts (# Correct), the number of incorrect facts (# Incorrect), the number of total facts (# Total) and the proportion of correct facts out of the total number of extracted facts (% Correct). More details of the evaluation are in Appendix C.

4.5 Implementation details

We employ Llama2-base models of different sizes (7b and 13b) as our backbone models for training. We adopt the Alpaca template (Taori et al., 2023) for training and inference. The OpenAI model used for Extract(\cdot), Rewrite(\cdot) and Revise(\cdot) is *gpt-3.5-turbo*. More details of the implementation are in Appendix D.

5 Experimental results and analysis

To answer our research questions, we conduct generic domain and medical domain QA experiments, fine-grained facts evaluation, and ablation studies. In addition, we conducted a case study to gain further understanding of the effectiveness of KnowTuning.

5.1 Main results (RQ1)

Automatic evaluation. Table 1 presents the reference-based GPT-4 evaluation results for both generic and medical domain QA datasets. Across all metrics, KnowTuning outperforms the baseline models in these domains. Based on the results, we

Method	Dataset	Completeness			Factuality			Logicity			Avg. gap
		Win	Tie	Lose	Win	Tie	Lose	Win	Tie	Lose	
Backbone Language Model: Llama2-7b-base											
KnowTuning vs FactTune	Dolly	61.00*	12.00	27.00	28.00*	58.50	13.50	33.50*	50.00	16.50	+21.83
KnowTuning vs FactTune	MedQuAD	73.00*	9.00	18.00	40.00*	43.00	17.00	45.50*	36.00	18.50	+35.00
Backbone Language Model: Llama2-13b-base											
KnowTuning vs FactTune	Dolly	58.00*	11.00	31.00	32.50*	66.50	11.00	35.00*	53.00	12.00	+23.83
KnowTuning vs FactTune	MedQuAD	78.00*	6.50	15.50	43.00*	45.50	11.50	39.00*	45.50	15.50	+39.17

Table 2: Human evaluation results on generic domain and medical domain QA datasets. The scores marked with * mean KnowTuning surpass FactTune significantly with p -value < 0.05 (sign. test).

Method	Dolly				MedQuAD			
	# Correct ↑	# Incorrect ↓	# Total ↑	% Correct ↑	# Correct ↑	# Incorrect ↓	# Total ↑	% Correct ↑
Backbone Language Model: Llama2-7b-base								
Base	6.15	3.62	9.77	62.94	6.54	3.42	9.96	65.66
SFT	7.77	1.85	9.62	80.77	16.11	1.73	17.84	90.30
RLAIF	11.23	2.10	13.33	84.25	10.86	0.95	11.81	91.95
FactTune	11.25	1.92	13.17	85.42	12.83	0.83	13.66	93.92
KnowTuning	14.40	2.36	16.76	85.89	18.04	0.98	19.02	94.87
Backbone Language Model: Llama2-13b-base								
Base	9.57	4.28	13.85	69.11	7.96	3.50	11.46	69.46
SFT	9.96	2.21	12.17	81.84	16.82	1.66	18.48	91.02
RLAIF	10.72	2.16	12.88	83.26	13.01	1.16	14.17	91.81
FactTune	12.73	2.12	14.85	85.72	13.02	1.01	14.03	92.80
KnowTuning	15.44	2.20	17.64	87.54	19.01	1.11	20.12	94.48

Table 3: Fine-grained facts evaluation on generic and medical QA. The best performance is highlighted in **bold**.

have two main observations:

- **KnowTuning consistently surpasses baselines in terms of completeness, factuality and logicity, across generic and domain-specific QA datasets.** Compared with Base and SFT, KnowTuning focuses on improving fine-grained and coarse-grained knowledge awareness of LLMs, which significantly improves the performance. Compared with RLAIF and FactTune, KnowTuning is more effective in improving the performance of LLMs on complex knowledge-intensive QA in multiple aspects. The reason is that RLAIF improves the performance by calculating overall helpfulness scores and FactTune focuses on improving the factuality, they ignore improving the knowledge awareness of LLMs in multiple essential aspects simultaneously.
- **KnowTuning demonstrates effectiveness on LLMs across different sizes.** We observe that KnowTuning consistently improves the performance of QA tasks on different scales (7b and 13B) LLMs. This finding aligns with [Bian et al. \(2024\)](#) and [Mecklenburg et al. \(2024\)](#): LLMs learn a lot of generic knowledge during the pre-training stage but still need to learn downstream domain knowledge and explore how to effec-

tively leverage knowledge for solving knowledge-intensive QA tasks.

Human evaluation. Human evaluations are crucial for accurately assessing the quality of answers. As shown in Table 2, to facilitate human annotation processes, we focus on comparing KnowTuning with the state-of-art baseline FactTune:

- Our findings indicate that KnowTuning consistently surpasses FactTune in terms of completeness, factuality, and logicity performance across various sizes of LLMs under human evaluation.
- KnowTuning demonstrates superior performance over QA in both generic and medical domain QA evaluated by human, in terms of completeness, factuality, and logicity.

5.2 Fine-grained facts evaluation (RQ2)

To evaluate the ability of methods to generate correct facts at the fine-grained level, we conduct fine-grained facts evaluation experiments. Based on the results in Table 3, we have two main observations:

- **Knowtuning generates answers with a higher proportion of correct facts across various sizes.** Compared to baselines, KnowTuning can generate more facts with less factual error rate across

Method	Completeness			Factuality			Logicity			Avg. gap
	Win	Tie	Lose	Win	Tie	Lose	Win	Tie	Lose	
-KA vs KnowTuning	32.50	20.00	47.50	16.00	57.50	26.50	12.50	61.50	26.00	-13.00
-KCC vs KnowTuning	18.50	31.00	50.50	11.00	72.50	16.50	10.50	61.50	28.00	-18.33
-KFC vs KnowTuning	23.00	28.50	48.50	8.50	70.50	21.00	12.00	60.50	27.50	-17.83
-KLC vs KnowTuning	25.50	27.50	47.00	12.00	73.00	15.00	9.50	60.00	30.50	-15.17
-KC vs KnowTuning	11.50	6.00	82.50	16.00	52.00	32.00	15.50	40.50	44.00	-38.50

Table 4: Ablation study evaluated by GPT-4 on the generic QA dataset. The backbone model is Llama2-7b-base.

different sizes of LLMs. Although RLAIIF and FactTune improve the proportion of correct facts, they ignore fine-grained knowledge augmentation and coarse-grained knowledge completeness awareness. Note that even though FactTune generates fewer incorrect facts, KnowTuning outperforms FactTune on the more critical metric of the percentage of correct facts.

- **KnowTuning generates larger amounts of correct facts across generic and domain-specific QA datasets.** Compared to SFT, we observe that KnowTuning consistently generates more correct facts across generic and domain-specific QA datasets. However, in the specific medical domain QA, RLAIIF and FactTune generate fewer correct facts than SFT. This is because LLMs learn a large amount of generic knowledge during the pre-training stage, yet still lack domain-specific knowledge for downstream tasks (Mecklenburg et al., 2024). This underscores the necessity for enhancing fine-grained knowledge awareness in domain-specific, knowledge-intensive QA tasks, as well as the need to improve coarse-grained knowledge awareness across key aspects of completeness, factuality, and logicity.

5.3 Ablation studies (RQ3)

In Table 4, we compare KnowTuning with several ablative variants. The variants are as follows: (i) **-KA**: we remove the fine-grained knowledge augmentation. (ii) **-KCC**: we remove knowledge completeness comparison set. (iii) **-KFC**: we remove knowledge factuality comparison set. (iv) **-KLC**: we remove knowledge logicity comparison set. (v) **-KC**: we remove all coarse-grained knowledge comparison sets. Our findings are as follows:

- **Removing the fine-grained knowledge augmentation.** We observe that removing fine-grained knowledge augmentation (-KA) decreases the performance of all three aspects. This indicates that fine-grained knowledge augmentation is effective for improving fine-grained

knowledge awareness of LLMs.

- **Removing the coarse-grained knowledge comparison.** The absence of coarse-grained knowledge comparisons results in substantial performance degradation in knowledge-intensive QA tasks. Specifically, removing the knowledge completeness comparison (-KCC) adversely affects completeness, the elimination of the knowledge factuality comparison (-KFC) undermines factuality, and the removal of the knowledge logicity comparison (-KLC) diminishes logicity. Although deleting and revising atomic knowledge can impact logicity, shuffling has been found more effective in improving coarse-grained logicity for LLMs. Furthermore, removing all coarse-grained knowledge comparison sets (-KC) results in a significant drop in performance across all aspects of the knowledge-intensive QA task.

5.4 Case study

We conduct several case studies and find that KnowTuning is more effective at generating complete, factual and logical answers than baselines across various sizes of LLMs. More details of our case study results are in Appendix E.

6 Conclusions

In this paper, we focus on improving the knowledge awareness of LLMs via fine-tuning for complex knowledge-intensive tasks. We have proposed KnowTuning to fine-tune LLMs through fine-grained knowledge augmentation and coarse-grained knowledge comparison stages. We have conducted comprehensive experiments on generic and medical domain QA datasets, demonstrating the effectiveness of KnowTuning through automatic and human evaluations, across various sizes of LLMs. Moreover, KnowTuning generates more facts with less factual error rate under fine-grained facts evaluation. Our code and dataset are available at <https://anonymous.4open.science/r/KnowTuning-345D>.

598 Limitations

599 In this study, KnowTuning is mainly aimed at
600 generic and medical knowledge-intensive tasks, we
601 plan to adopt KnowTuning to other tasks such as
602 legal domain QA (Zhong et al., 2020) and math-
603 ematical reasoning (Luo et al., 2023). Moreover,
604 our efforts have been concentrated on enhancing
605 the knowledge awareness of LLMs during the fine-
606 tuning stage. Future studies will aim to explore
607 improving knowledge awareness of LLMs in the
608 pre-training stage (Rosset et al., 2020).

609 Ethics Statement

610 KnowTuning mainly focuses on completeness, fac-
611 tuality, and logicity, but not social bias or the po-
612 tential for generating harmful or toxic content (He-
613 witt et al., 2024). We plan to adopt our method
614 to reduce social bias and harmful content at fine-
615 grained and coarse-grained levels in future work.

616 References

617 Asma Ben Abacha and Dina Demner-Fushman. 2019. A
618 question-entailment approach to question answering.
619 *BMC Bioinform.*, 20(1):511:1–511:23.

620 Tal August, Katharina Reinecke, and Noah A. Smith.
621 2022. Generating scientific definitions with control-
622 lable complexity. In *Proceedings of ACL*, pages
623 8298–8317.

624 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
625 Amanda Askell, Jackson Kernion, Andy Jones, Anna
626 Chen, Anna Goldie, Azalia Mirhoseini, Cameron
627 McKinnon, Carol Chen, Catherine Olsson, Christo-
628 pher Olah, Danny Hernandez, Dawn Drain, Deep
629 Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez,
630 Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua
631 Landau, Kamal Ndousse, Kamile Lukosiute, Liane
632 Lovitt, Michael Sellitto, Nelson Elhage, Nicholas
633 Schiefer, Noemí Mercado, Nova DasSarma, Robert
634 Lasenby, Robin Larson, Sam Ringer, Scott John-
635 ston, Shauna Kravec, Sheer El Showk, Stanislav Fort,
636 Tamera Lanham, Timothy Telleen-Lawton, Tom Con-
637 erly, Tom Henighan, Tristan Hume, Samuel R. Bow-
638 man, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,
639 Nicholas Joseph, Sam McCandlish, Tom Brown, and
640 Jared Kaplan. 2022. Constitutional AI: harmless-
641 ness from AI feedback. *CoRR*, abs/2212.08073.

642 Yuyang Bai, Shangbin Feng, Vidhisha Balachandran,
643 Zhaoxuan Tan, Shiqi Lou, Tianxing He, and Yulia
644 Tsvetkov. 2023. Kgquiz: Evaluating the generaliza-
645 tion of encoded knowledge in large language models.
646 *CoRR*, abs/2310.09725.

647 Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie
648 Lu, and Ben He. 2024. Chatgpt is a knowledge-
649 able but inexperienced solver: An investigation of

commonsense problem in large language models. In
Proceedings of COLING.

650
651

652 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
653 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
654 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
655 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
656 Gretchen Krueger, Tom Henighan, Rewon Child,
657 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
658 Clemens Winter, Christopher Hesse, Mark Chen, Eric
659 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
660 Jack Clark, Christopher Berner, Sam McCandlish,
661 Alec Radford, Ilya Sutskever, and Dario Amodei.
662 2020. Language models are few-shot learners. In
663 *Proceedings of NeurIPS*.

664 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
665 Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi,
666 Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang,
667 Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie.
668 2023. A survey on evaluation of large language mod-
669 els. *CoRR*, abs/2307.03109.

670 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa
671 Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniv-
672 asan, Tianyi Zhou, Heng Huang, and Hongxia Jin.
673 2023a. Alpapasus: Training a better Alpaca with
674 fewer data. *CoRR*, abs/2307.08701.

675 Shiqi Chen, Yiran Zhao, Jinghan Zhang, I-Chun Chern,
676 Siyang Gao, Pengfei Liu, and Junxian He. 2023b.
677 FELM: benchmarking factuality evaluation of large
678 language models. *CoRR*, abs/2310.00741.

679 Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Junchen
680 Wan, Fuzheng Zhang, Di Zhang, and Ji-Rong Wen.
681 2024. Improving large language models via fine-
682 grained reinforcement learning with minimum edit-
683 ing constraint. *CoRR*, abs/2401.06081.

684 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
685 Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang,
686 Mostafa Dehghani, Siddhartha Brahma, Albert Web-
687 son, Shixiang Shane Gu, Zhuyun Dai, Mirac Suz-
688 gun, Xinyun Chen, Aakanksha Chowdhery, Sharan
689 Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao,
690 Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav
691 Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam
692 Roberts, Denny Zhou, Quoc V. Le, and Jason Wei.
693 2022. Scaling instruction-finetuned language models.
694 *CoRR*, abs/2210.11416.

695 Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie,
696 Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell,
697 Matei Zaharia, and Reynold Xin. 2023. Free dolly:
698 Introducing the world’s first truly open instruction-
699 tuned llm.

700 Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan,
701 Shizhe Diao, Jipeng Zhang, Kashun Shum, and
702 Tong Zhang. 2023. RAFT: reward ranked finetuning
703 for generative foundation model alignment. *CoRR*,
704 abs/2304.06767.

705 Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun
706 Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao

707	Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. Loramoe: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment . <i>CoRR</i> , abs/2312.09979.	763
708		764
709		765
710		766
711		767
712	Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpaca-Farm: A simulation framework for methods that learn from human feedback . <i>CoRR</i> , abs/2305.14387.	768
713		769
714		770
715		771
716		772
717	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. GPTScore: Evaluate as you desire . <i>CoRR</i> , abs/2302.04166.	773
718		774
719		775
720	Jian Guan, Xiaoxi Mao, Changjie Fan, Zitao Liu, Wenbiao Ding, and Minlie Huang. 2021. Long text generation by modeling sentence-level and discourse-level coherence . In <i>Proceedings of ACL</i> , pages 6379–6393.	776
721		777
722		778
723		779
724		780
725	Geyang Guo, Ranchi Zhao, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Beyond imitation: Leveraging fine-grained quality signals for alignment . <i>CoRR</i> , abs/2311.04072.	781
726		782
727		783
728		784
729	John Hewitt, Sarah Chen, Lanruo Lora Xie, Edward Adams, Percy Liang, and Christopher D Manning. 2024. Model editing with canonical examples . <i>arXiv preprint arXiv:2402.06155</i> .	785
730		786
731		787
732		788
733	Hiyouga. 2023. Llama factory . https://github.com/hiyouga/LLaMA-Factory .	789
734		790
735	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models . In <i>Proceedings of ICLR</i> .	791
736		792
737		793
738		794
739	Wenyue Hua, Jiang Guo, Mingwen Dong, Henghui Zhu, Patrick Ng, and Zhiguo Wang. 2024. Propagation and pitfalls: Reasoning-based assessment of knowledge editing through counterfactual tasks . <i>CoRR</i> , abs/2401.17585.	795
740		796
741		797
742		798
743		799
744	Neel Jain, Ping-yeh Chiang, Yuxin Wen, John Kirchenbauer, Hong-Min Chu, Gowthami Somepalli, Brian R. Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Neptune: Noisy embeddings improve instruction finetuning . <i>CoRR</i> , abs/2310.05914.	800
745		801
746		802
747		803
748		804
749		805
750		806
751	Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. Towards mitigating hallucination in large language models via self-reflection . <i>CoRR</i> , abs/2310.06271.	807
752		808
753		809
754		810
755	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2023. Mistral 7b . <i>CoRR</i> , abs/2310.06825.	811
756		812
757		813
758		814
759		815
760		816
761		817
762		818
	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, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. 2022. Language models (mostly) know what they know . <i>CoRR</i> , abs/2207.05221.	763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
	Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. 2023. Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks . <i>CoRR</i> , abs/2305.18395.	776
		777
		778
		779
		780
	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive NLP . <i>CoRR</i> , abs/2212.14024.	781
		782
		783
		784
		785
	Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo Kim, and Jaewoo Kang. 2020. Look at the first sentence: Position bias in question answering . In <i>Proceedings of EMNLP</i> , pages 1109–1121.	786
		787
		788
		789
	Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering . In <i>Proceedings of NAACL-HLT</i> , pages 4940–4957.	790
		791
		792
	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbone, and Abhinav Rastogi. 2023. RLAIF: scaling reinforcement learning from human feedback with AI feedback . <i>CoRR</i> , abs/2309.00267.	793
		794
		795
		796
		797
	Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale Fung, Mohammad Shoeybi, and Bryan Catanzaro. 2022. Factuality enhanced language models for open-ended text generation . In <i>Proceedings of NeurIPS</i> .	798
		799
		800
		801
		802
	Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2023. From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning . <i>CoRR</i> , abs/2308.12032.	803
		804
		805
		806
		807
	Jiacheng Liu, Skyler Hallinan, Ximing Lu, Pengfei He, Sean Welleck, Hannaneh Hajishirzi, and Yejin Choi. 2022a. Rainier: Reinforced knowledge introspector for commonsense question answering . In <i>Proceedings of EMNLP</i> , pages 8938–8958.	808
		809
		810
		811
		812
	Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2022b. Generated knowledge prompting for commonsense reasoning . In <i>Proceedings of ACL</i> , pages 3154–3169.	813
		814
		815
		816
		817
		818

818	Jiacheng Liu, Ramakanth Pasunuru, Hannaneh Hajishirzi, Yejin Choi, and Asli Celikyilmaz. 2023a. Crystal: Introspective reasoners reinforced with self-feedback . In <i>Proceedings of EMNLP</i> , pages 11557–11572.	874
819		875
820		876
821		877
822		878
823	Yixin Liu, Alexander R. Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023b. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation . In <i>Proceedings of ACL</i> , pages 4140–4170. Association for Computational Linguistics.	879
824		880
825		881
826		882
827		
828		
829		
830		
831	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. The flan collection: Designing data and methods for effective instruction tuning . In <i>Proceedings of ICML</i> , volume 202, pages 22631–22648.	883
832		884
833		885
834		886
835		887
836		888
837	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>Proceedings of ICLR</i> .	889
838		890
839	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct . <i>arXiv preprint arXiv:2308.09583</i> .	891
840		
841		
842		
843		
844		
845	Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods . https://github.com/huggingface/peft .	892
846		893
847		894
848		895
849		896
850	Nick Mecklenburg, Yiyou Lin, Xiaoxiao Li, Daniel Holstein, Leonardo Nunes, Sara Malvar, Bruno Silva, Ranveer Chandra, Vijay Aski, Pavan Kumar Reddy Yannam, et al. 2024. Injecting new knowledge into large language models via supervised fine-tuning . <i>arXiv preprint arXiv:2404.00213</i> .	897
851		898
852		899
853		900
854		901
855		902
856	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation . In <i>Proceedings of EMNLP</i> , pages 12076–12100.	903
857		904
858		905
859		906
860		907
861		908
862	Aditi Mishra, Sajjadur Rahman, Hannah Kim, Kushan Mitra, and Estevam Hruschka. 2023. Characterizing large language models as rationalizers of knowledge-intensive tasks . <i>CoRR</i> , abs/2311.05085.	909
863		
864		
865		
866	Fedor Moiseev, Zhe Dong, Enrique Alfonseca, and Martin Jaggi. 2022. SKILL: Structured knowledge infusion for large language models . In <i>Proceedings of NAACL</i> , pages 1581–1588.	910
867		911
868		912
869		913
870	Ani Nenkova and Rebecca J. Passonneau. 2004. Evaluating content selection in summarization: The pyramid method . In <i>Proceedings of HLT-NAACL</i> , pages 145–152.	914
871		915
872		916
873		
	Minh Nguyen, Kishan K. C., Toan Nguyen, Ankit Chadha, and Thuy Vu. 2023. Efficient fine-tuning large language models for knowledge-aware response planning . In <i>Machine Learning and Knowledge Discovery in Databases: Research Track - European Conference, ECML PKDD 2023, Turin, Italy, September 18-22, 2023, Proceedings, Part II</i> , volume 14170 of <i>Lecture Notes in Computer Science</i> , pages 593–611.	917
		918
		919
		920
	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 <i>Proceedings of NeurIPS</i> .	921
		922
		923
		924
		925
		926
		927
		928
		929
		930
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	
	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with GPT-4 . <i>CoRR</i> , abs/2304.03277.	
	Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? In <i>Proceedings of EMNLP</i> , pages 1339–1384.	
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model . <i>CoRR</i> , abs/2305.18290.	
	Corby Rosset, Chenyan Xiong, Minh Phan, Xia Song, Paul N. Bennett, and Saurabh Tiwary. 2020. Knowledge-aware language model pretraining . <i>CoRR</i> , abs/2007.00655.	
	Lin CY ROUGE. 2004. A package for automatic evaluation of summaries . In <i>Proceedings of Workshop on Text Summarization of ACL</i> , volume 5.	
	Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross J. Anderson. 2023. The curse of recursion: Training on generated data makes models forget . <i>CoRR</i> , abs/2305.17493.	
	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Schärli, Aakanksha Chowdhery, Philip Andrew Mansfield, Blaise Agüera y Arcas, Dale R. Webster, Gregory S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle K. Barral, Christopher Semurs, Alan Karthikesalingam, and Vivek	

931	Natarajan. 2022. Large language models encode clinical knowledge . <i>CoRR</i> , abs/2212.13138.	language models: Knowledge, retrieval and domain-specificity. <i>CoRR</i> , abs/2310.07521.	988
932			989
933	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2023. Preference ranking optimization for human alignment . <i>CoRR</i> , abs/2306.17492.	Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023c. Knowledge-driven CoT: Exploring faithful reasoning in LLMs for knowledge-intensive question answering . <i>CoRR</i> , abs/2308.13259.	990
934			991
935			992
936			993
937	Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. 2022. Read before generate! Faithful long form question answering with machine reading . In <i>Findings of ACL</i> , pages 744–756.	Peiyi Wang, Lei Li, Liang Chen, Feifan Song, Binghuai Lin, Yunbo Cao, Tianyu Liu, and Zhifang Sui. 2023d. Making large language models better reasoners with alignment . <i>CoRR</i> , abs/2309.02144.	994
938			995
939			996
940			997
941			998
942	Weiwei Sun, Zhengliang Shi, Shen Gao, Pengjie Ren, Maarten de Rijke, and Zhaochun Ren. 2023. Contrastive learning reduces hallucination in conversations . In <i>Proceedings of AAAI</i> , pages 13618–13626.	Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023e. Large language models are not fair evaluators . <i>CoRR</i> , abs/2305.17926.	999
943			1000
944			1001
945			1002
946	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model . https://github.com/tatsu-lab/stanford_alpaca .	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023f. Self-instruct: Aligning language models with self-generated instructions . In <i>Proceedings of ACL</i> , pages 13484–13508.	1003
947			1004
948			1005
949			1006
950			1007
951	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2023. Fine-tuning language models for factuality . In <i>Proceedings of NeurIPS Workshop on Instruction Tuning and Instruction Following</i> .	Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ NLP tasks . In <i>Proceedings of EMNLP</i> , pages 5085–5109.	1008
952			1009
953			1010
954			1011
955			1012
956	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models . <i>CoRR</i> , abs/2307.09288.	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners . In <i>Proceedings of ICLR</i> .	1013
957			1014
958			1015
959			1016
960			1017
961			1018
962			1019
963			1020
964			1021
965			1022
966			1023
967			1024
968			1025
969			1026
970			1027
971			1028
972			1029
973			1030
974			1031
975			1032
976			1033
977			1034
978			1035
979	Cunxiang Wang, Sirui Cheng, Zhikun Xu, Bowen Ding, Yidong Wang, and Yue Zhang. 2023a. Evaluating open question answering evaluation . <i>CoRR</i> , abs/2305.12421.	Fangyuan Xu, Yixiao Song, Mohit Iyyer, and Eunsol Choi. 2023a. A critical evaluation of evaluations for long-form question answering . In <i>Proceedings of ACL</i> , pages 3225–3245.	1036
980			1037
981			1038
982			1039
983	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Jiayang Cheng, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang, Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang, and Yue Zhang. 2023b. Survey on factuality in large	Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-Seng Chua. 2023b. Search-in-the-chain: Towards the accurate, credible and traceable content generation for complex knowledge-intensive tasks . <i>CoRR</i> , abs/2304.14732.	1040
984			1041
985			1042
986			1043
987			1044

1042	Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao,	pairs specifically categorized under "open_qa"	1097
1043	Daniel Zhang-li, Xin Lv, Hao Peng, Zijun Yao, Xi-	and "general_qa" for our dataset. We filter 4,000	1098
1044	aohan Zhang, Hanming Li, Chunyang Li, Zheyuan	QA pairs for training, 200 QA pairs for valida-	1099
1045	Zhang, Yushi Bai, Yantao Liu, Amy Xin, Nianyi Lin,	tion, and 200 QA pairs for testing.	1100
1046	Kaifeng Yun, Linlu Gong, Jianhui Chen, Zhili Wu,	• MedQuAD (Abacha and Demner-Fushman,	1101
1047	Yunjia Qi, Weikai Li, Yong Guan, Kaisheng Zeng,	2019): The dataset covers 37 different question	1102
1048	Ji Qi, Hailong Jin, Jinxin Liu, Yu Gu, Yuan Yao,	types. In this paper, following (August et al.,	1103
1049	Ning Ding, Lei Hou, Zhiyuan Liu, Bin Xu, Jie Tang,	2022), we filter QA pairs of the category "Infor-	1104
1050	and Juanzi Li. 2023a. Kola: Carefully benchmarking	information" for giving definitions and information	1105
1051	world knowledge of large language models. <i>CoRR</i> ,	about medical terms. We filter 4000 QA pairs	1106
1052	abs/2306.09296.	for training, 200 QA pairs for validation and 200	1107
1053	Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu,	QA pairs for testing.	1108
1054	Mingxuan Ju, Soumya Sanyal, Chenguang Zhu,	B Details of Baselines	1109
1055	Michael Zeng, and Meng Jiang. 2023b. Gener-	• Base: We adopt the Alpaca template (Taori et al.,	1110
1056	ate rather than retrieve: Large language models are	2023) for testing the Llama2-base model (Tou-	1111
1057	strong context generators. In <i>Proceedings of ICLR</i> .	vron et al., 2023) under zero-shot setting.	1112
1058	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang,	• SFT: We follow standard vanilla fine-tuning loss	1113
1059	Songfang Huang, and Fei Huang. 2023. RRHF: rank	in Eq. 1 to train LLMs on original QA datasets.	1114
1060	responses to align language models with human feed-	• RLAIF (Bai et al., 2022; Lee et al., 2023): We	1115
1061	back without tears. <i>CoRR</i> , abs/2304.05302.	leverage <i>gpt-3.5-turbo</i> to annotate overall help-	1116
1062	Shiyue Zhang and Mohit Bansal. 2021. Finding a bal-	fulness scores and construct generic helpfulness	1117
1063	anced degree of automation for summary evaluation.	comparison sets. We adopt DPO (Rafailov et al.,	1118
1064	In <i>Proceedings of EMNLP</i> , pages 6617–6632.	2023) for generic helpfulness comparison sets	1119
1065	Yichi Zhang, Zhuo Chen, Yin Fang, Lei Cheng, Yanxi	optimization.	1120
1066	Lu, Fangming Li, Wen Zhang, and Huajun Chen.	• FactTune (Tian et al., 2023): We follow Min et al.	1121
1067	2023a. Knowledgeable preference alignment for	(2023) to first break each candidate answers into	1122
1068	llms in domain-specific question answering. <i>CoRR</i> ,	individual facts, and prompt LLMs to measure	1123
1069	abs/2311.06503.	the correctness of each fact based on the golden	1124
1070	Yue Zhang, Ming Zhang, Haipeng Yuan, Shichun Liu,	answer as a reference. ¹ Then, we construct fac-	1125
1071	Yongyao Shi, Tao Gui, Qi Zhang, and Xuanjing	tuality comparison sets by the percentage of cor-	1126
1072	Huang. 2023b. Llmeval: A preliminary study on	rect facts. Finally, we adopt DPO (Rafailov et al.,	1127
1073	how to evaluate large language models. <i>CoRR</i> ,	2023) for factuality comparison sets optimiza-	1128
1074	abs/2312.07398.	tion.	1129
1075	Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman,	C Details of Evaluation	1130
1076	Mohammad Saleh, and Peter J. Liu. 2023. Slic-hf:	C.1 GPT-4 Evaluation	1131
1077	Sequence likelihood calibration with human feed-	This section provides specifics of the GPT-4 prompt	1132
1078	back. <i>CoRR</i> , abs/2305.10425.	utilized for reference-based evaluation, employing	1133
1079	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	<i>gpt4-turbo</i> . Figure 3 illustrates the adapted prompt	1134
1080	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	from Zheng et al. (2024), aimed at assessing the	1135
1081	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	completeness, factuality, and logicity of answers.	1136
1082	Joseph E. Gonzalez, and Ion Stoica. 2024. Judg-	To avoid positional bias (Ko et al., 2020; Wang	1137
1083	ing llm-as-a-judge with mt-bench and chatbot arena.	et al., 2023e), we evaluate each answer in both	1138
1084	<i>Proceedings of NeurIPS</i> , 36.	positions during two separate runs.	1139
1085	Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang	C.2 Human Evaluation	1140
1086	Zhang, Zhiyuan Liu, and Maosong Sun. 2020. Jec-	For the human evaluation, we hired people with	1141
1087	qa: a legal-domain question answering dataset. In	undergraduate degrees and undergraduate medical	1142
1088	<i>Proceedings of AACL</i> , volume 34, pages 9701–9708.	degrees to annotate generic QA and medical QA	1143
1089	Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and		
1090	Dacheng Tao. 2023. Can chatgpt understand too? A		
1091	comparative study on chatgpt and fine-tuned BERT.		
1092	<i>CoRR</i> , abs/2302.10198.		
1093	Appendix		
1094	A Details of Datasets		
1095	• Dolly (Conover et al., 2023): Given our focus on		
1096	open-ended generic domain QA, we selected QA		

¹<https://github.com/shmsw25/FActScore>

[System prompt]
 You are a helpful and precise assistant for checking the quality of the answer.

[User prompt]
 [Question]
 {question}

[The Start of Reference Answer]
 {answer_ref}
 [The End of Reference Answer]

[The Start of Assistant 1's response]
 {answer_a}
 [The End of Assistant 1's response]

[The Start of Assistant 2's response]
 {answer_b}
 [The End of Assistant 2's response]

We would like to request your feedback on the performance of two AI assistants in response to the user question displayed above. Based the reference answer, you should rate the Knowledge Completeness, Knowledge Factuality and Knowledge Logicality of their responses. Each aspect of each assistant receives an score on a scale of 1 to 10, where a higher score indicates better performance. Please generate Knowledge Completeness, Knowledge Factuality and Knowledge Logicality scores for each assistant in order.

Please generate the scores in order and following format.
 {'Knowledge Completeness':value,'Knowledge Factuality':value,'Knowledge Logicality':value}

Please first output two lines containing values indicating the Knowledge Completeness, Knowledge Factuality and Knowledge Logicality scores for Assistant 1 and 2, respectively. In the subsequent line, please provide a comprehensive explanation of your evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not affect your judgment.

Figure 3: Prompts for GPT-4 evaluation.

You'll be presented with a series of questions. For each question, two answers and a golden answer will be provided. Your task is to read both answers carefully and decide which one you believe is better.

When judging, consider:

- Completeness: It examines whether the answers provide comprehensive and sufficient knowledge relevant to the questions.
- Factuality: It examines whether the knowledge in the answers is factually correct
- Logicality: it examines whether the knowledge in the answers is logically rigorous and structured.

Question:
 {Q}

Golden Answer:
 {A0}

Answer A:
 {A1}

Answer B:
 {A2}

Based on the golden answer, comparing these two answers, in terms of completeness, factuality and logicality, respectively.
 Give the win-tie-lose of Answer A compared to Answer B in each of the three aspects.

Figure 4: Instructions for human evaluation.

test sets, respectively, to ensure the trustworthiness of the human evaluations, and we allowed the human evaluators to access Wikipedia to further validate the knowledge during the evaluation process. Instructions for human evaluation are depicted in Figure 4.

C.3 Fine-grained facts evaluation

Following Min et al. (2023), we first break candidate answers into individual facts, and use *gpt-3.5-turbo* to measure the correctness of each fact based on the golden answer as a reference.¹

D Details of Implementation

D.1 Prompts for Extracting, Rewriting, and Revising

Details for the prompts used in `Extract(·)`, `Rewrite(·)`, and `Revise(·)` are provided. Figures 5, 6, 7 and 8 display the prompts for extracting atomic knowledge, rewriting fine-grained questions, rewriting fine-grained answers, and revising atomic knowledge into nonfactual knowledge, respectively.

D.2 Reliability of atomic knowledge extraction

To evaluate the reliability of atomic knowledge extraction, we first sample 50 instances of genericQA dataset Dolly. We manually checked these data and find that only 3 instances required further separation or merging of atomic facts, illustrating the reliability of extracting atomic facts using *gpt3.5-turbo*.

D.3 Training

During the training phase, the AdamW optimizer (Loshchilov and Hutter, 2019) is utilized with initial learning rates of $1 \cdot 10^{-4}$ for SFT and $5 \cdot 10^{-6}$ for DPO. The batch sizes for SFT and DPO are set to 32 and 16, respectively, with SFT undergoing 3 epochs of training and DPO 1 epoch. The filtering and deleting percentages, α and β , are both fixed at 0.5. We determine the hyperparameters through pilot experiments. Training leverages PEFT (Mangrulkar et al., 2022), LLaMA-Factory (Hyouga, 2023) and LoRA (Hu et al., 2022).

D.4 Cost Analysis

The cost of KnowTuning is lower than that of the baseline methods RLAIIF and FactTune. Specifically, in the generic domain QA dataset Dolly, the

costs are as follows: KnowTuning is \$8.45, RLAIIF is \$9.94, and FactTune is \$10.53. This cost difference arises because RLAIIF necessitates pairwise comparisons for assessing the overall helpfulness of all candidate answers, while FactTune requires a detailed factuality evaluation for each fact across all candidate answers, thereby increasing their dataset comparison construction costs.

E Details of Case Study

As illustrated in Figures 9 and 10, the case studies evaluate answers generated by four methods: SFT, RLAIIF, FactTune, and KnowTuning across various sizes. Our findings indicate that KnowTuning excels at producing answers that are more complete, factual, and logical across various sizes of LLMs, as detailed below:

- As shown in Figure 9 for the case study based on backbone Llama2-7b-base, KnowTuning generates more complete and logical answers compared to all baselines. Although RLAIIF produces more knowledge compared to SFT, it results in fewer logical answers because it does not explicitly focus on logic optimization. FactTune, on the other hand, focuses on improving the percentage of factualness and performs poorly in terms of answer completeness and logic. This illustrates the need for multiple aspects of coarse-grained knowledge awareness.
- As shown in Figure 10 for the case study based on backbone Llama2-13b-base, KnowTuning generates content that is more informative and factual, and the logic between the knowledge is more logical. Although RLAIIF generates multiple aspects of knowledge, it does not provide fine-grained knowledge in the answer. FactTune generates detailed information such as Canada’s domestic population and GDP, but it provides factually incorrect information. This further underscores the critical need for enhanced fine-grained knowledge awareness.

Please breakdown the following sentence into independent facts: He made his acting debut in the film *The Moon is the Sun's Dream* (1992), and continued to appear in small and supporting roles throughout the 1990s.

- He made his acting debut in the film.
- He made his acting debut in *The Moon is the Sun's Dream*.
- *The Moon is the Sun's Dream* is a film.
- *The Moon is the Sun's Dream* was released in 1992.
- After his acting debut, he appeared in small and supporting roles.
- After his acting debut, he appeared in small and supporting roles throughout the 1990s.

Please breakdown the following sentence into independent facts: He is also a successful producer and engineer, having worked with a wide variety of artists, including Willie Nelson, Tim McGraw, and Taylor Swift.

- He is successful.
- He is a producer.
- He is an engineer.
- He has worked with a wide variety of artists. - Willie Nelson is an artist.
- He has worked with Willie Nelson.
- Tim McGraw is an artist.
- He has worked with Tim McGraw.
- Taylor Swift is an artist.
- He has worked with Taylor Swift.

Please breakdown the following sentence into independent facts: In 1963, Collins became one of the third group of astronauts selected by NASA and he served as the back-up Command Module Pilot for the Gemini 7 mission.

- Collins became an astronaut.
- Collins became one of the third group of astronauts.
- Collins became one of the third group of astronauts selected.
- Collins became one of the third group of astronauts selected by NASA.
- Collins became one of the third group of astronauts selected by NASA in 1963. - He served as the Command Module Pilot.
- He served as the back-up Command Module Pilot.
- He served as the Command Module Pilot for the Gemini 7 mission.

Please breakdown the following sentence into independent facts: In addition to his acting roles, Bateman has written and directed two short films and is currently in development on his feature debut.

- Bateman has acting roles.
- Bateman has written two short films.
- Bateman has directed two short films.
- Bateman has written and directed two short films.
- Bateman is currently in development on his feature debut.

Please breakdown the following sentence into independent facts: Michael Collins (born October 31, 1930) is a retired American astronaut and test pilot who was the Command Module Pilot for the Apollo 11 mission in 1969.

- Michael Collins was born on October 31, 1930.
- Michael Collins is retired.
- Michael Collins is an American.
- Michael Collins was an astronaut.
- Michael Collins was a test pilot.
- Michael Collins was the Command Module Pilot.
- Michael Collins was the Command Module Pilot for the Apollo 11 mission.
- Michael Collins was the Command Module Pilot for the Apollo 11 mission in 1969.

Please breakdown the following sentence into independent facts: He was an American composer, conductor, and musical director. - He was an American.

- He was a composer.
- He was a conductor.
- He was a musical director.

Please breakdown the following sentence into independent facts: She currently stars in the romantic comedy series, *Love and Destiny*, which premiered in 2019. - She currently stars in *Love and Destiny*.

- *Love and Destiny* is a romantic comedy series.
- *Love and Destiny* premiered in 2019.

Please breakdown the following sentence into independent facts: During his professional career, McCoy played for the Broncos, the San Diego Chargers, the Minnesota Vikings, and the Jacksonville Jaguars.

- McCoy played for the Broncos.
- McCoy played for the Broncos during his professional career.
- McCoy played for the San Diego Chargers.
- McCoy played for the San Diego Chargers during his professional career. - McCoy played for the Minnesota Vikings.
- McCoy played for the Minnesota Vikings during his professional career. - McCoy played for the Jacksonville Jaguars.
- McCoy played for the Jacksonville Jaguars during his professional career.

Please breakdown the following sentence into independent facts

Figure 5: Prompts for extracting atomic knowledge in the answer (Min et al., 2023).

[System prompt]

I want you to act as an Excellent Rewriter. Your objective is to rewrite a specific question that asks for knowledge of the relevant aspects of the given facts. Please read the example carefully and follow the format of the example to generate it.

[User prompt]

#Example#:

#Given Facts#:

- Sandworms are huge.
- Sandworms are aggressive.
- Sandworms live in the sand seas.

#Rewritten Question#:

- What is the size, aggressiveness, and habitat of sandworms?

#Example#:

#Given Facts#:

- A Series I-Bond helps protect from inflation.
- The inflation rate is determined by the treasury department.
- The inflation rate is adjusted twice a year.

#Rewritten Question#:

- In terms of inflation protection, how does a Series I-Bond function, who sets its inflation rate, and how often is this rate reviewed and adjusted?

#Example#:

#Given Facts#:

- An apple is produced by an apple tree.
- Apple trees are cultivated worldwide.

#Rewritten Question#:

- How is the apple produced by apple trees, and what is the scope of their cultivation globally?

You should rewrite the given question using the following rules:

You should try your best not to make the #Rewritten Question# become verbose.

#Rewritten Question# can only add 10 to 20 words into #Given Question#.

#Rewritten Question# should contain more specific relevant intentions to the #Given Facts#.

'#Given Question#', '#Rewritten Question#', 'given question', and 'rewritten question' are not allowed to appear in #Rewritten Question#.

#Given Facts#:

{difficult facts}

#Rewritten Question#:

Figure 6: Prompts for rewriting fine-grained questions.

[System prompt]

I want you to act as a helpful assistant. Your objective is to rewrite a high-quality answer to the given question based on the given facts.

[User prompt]

#Given Question#:

{fine-grained question}

#Given Facts#:

{difficult facts}

#Answer#:

Figure 7: Prompts for rewriting fine-grained answers.

[System prompt]

I want you to act as an Excellent Reviser. Your objective is to revise the given facts into incorrect facts. Please read the example carefully and follow the examples to generate it.

[User prompt]

#Example#

#Given Facts#:

- Sandworms are huge.
- Sandworms are aggressive.
- Sandworms live in the sand seas.

#Incorrect Facts#:

- Sandworms are tiny.
- Sandworms are timid.
- Sandworms live in the ocean.

#Example#

#Given Facts#:

- A Series I-Bond helps protect from inflation.
- The inflation rate is determined by the treasury department.
- The inflation rate is adjusted twice a year.

#Incorrect Facts#:

- A Series I-Bond exacerbates inflation.
- The inflation rate is determined by random selection.
- The inflation rate is adjusted once every decade.

#Example#

#Given Facts#:

- An apple is produced by an apple tree.
- Apple trees are cultivated worldwide.

#Incorrect Facts#:

- A pineapple is produced by an apple tree.
- Apple trees are only found in Antarctica

You should revise the given facts using the following rules:

The number of #Incorrect Facts# has to be the same as the #Given Facts#

#Given Facts#:

{atomic facts}

#Incorrect Facts#:

Figure 8: Prompts for revising atomic facts into incorrect facts.

Question:

Name different kinds of cheese

SFT:

There are many kinds of cheese. They are typically categorized by the region where they originate. For example, cheddar cheese originates from England, whereas Gouda cheese originates from the Netherlands. Other kinds of cheese include Swiss, Parmesan, mozzarella, Brie, Camembert, Roquefort, and blue cheese.

RLAIF:

There are many different kinds of cheese, each with its own unique flavor and texture. Some of the most popular types of cheese include cheddar, mozzarella, brie, gouda, and feta. Cheddar is a hard cheese that is often used in grilled cheese sandwiches and macaroni and cheese. Mozzarella is a soft, white cheese that is often used in pizza. Brie is a soft, creamy cheese that is often used in desserts. Gouda is a hard, yellow cheese that is often used in grilled cheese sandwiches. Feta is a soft, white cheese that is often used in salads.

FactTune:

There are many different kinds of cheese, including cheddar, brie, gouda, mozzarella, blue cheese, and feta. Each type of cheese has its own unique flavor, texture, and aroma.

KnowTuning:

There are many different kinds of cheese, each with its own unique flavor and texture. Some common types of cheese include:

1. Mozzarella: A soft, stretchy cheese that is commonly used in pizza and pasta dishes.
2. Cheddar: A sharp, tangy cheese that is often used in grilled cheese sandwiches and macaroni and cheese.
3. Gouda: A creamy, slightly sweet cheese that is popular in grilled cheese sandwiches and salads.
4. Brie: A soft, creamy cheese that is often served with fruit or crackers.
5. Camembert: A creamy, slightly runny cheese that is popular in French cuisine.
6. Feta: A tangy, salty cheese that is commonly used in Greek salads and as a topping for pasta dishes.
7. Parmesan: A hard, sharp cheese that is often grated over pasta dishes or used in sauces.
8. Cheddar: A hard, crumbly cheese that is often used in grilled cheese sandwiches and macaroni and cheese.
9. Gruyere: A nutty, slightly sweet cheese that is popular in Swiss dishes and as a topping for mashed potatoes.
10. Parmigiano-Reggiano: A hard, aged cheese that is popular in Italian cuisine and often grated over pasta dishes.
11. Emmental: A Swiss cheese known for its distinctive holes, often used in fondue and sandwiches.
12. Edam: A mild, semi-hard cheese that is popular in sandwiches and salads.
13. Gorgonzola: A blue cheese with a tangy flavor that is often served with fruit or crackers.
14. Fontina: A creamy, nutty cheese that is often used in pasta dishes and fondue.
15. Ricotta: A soft, creamy cheese that is commonly used in pasta dishes and as a filling for ravioli.

Figure 9: Case study for intuitive comparisons on generic QA dataset based on Llama2-7b-base.

Question:

What is Canada known for?

SFT:

Canada is known for its beautiful landscape, which is the result of a combination of natural geographic features and human-made efforts to preserve nature.

Canada is also known for its winters, which are long and cold in most regions. Canadians have adapted to these harsh winters and enjoy a variety of winter sports.

Canada is also known for its diverse population and its tolerance for different cultures and religions. This has led to a vibrant multicultural society that is open to new ideas and experiences.

RLAIF:

Canada is known for its diverse landscape, including the Rocky Mountains, the Canadian Shield, and the Great Lakes. Canada is also known for its natural resources, including oil, gas, and minerals. Canada is also known for its political stability and strong economy.

FactTune:

Canada is known for its natural beauty, including the Rocky Mountains, the Great Lakes, and the Canadian Rockies. Canada is also known for its cultural diversity, with a population of over 37 million people from a variety of backgrounds. Canada is also known for its strong economy, with a GDP of over \$1.8 trillion in 2020. Canada is also known for its high quality of life, with a high standard of living and a strong social safety net.

KnowTuning:

Canada is known for its stunning natural beauty, including vast forests, pristine lakes, and snow-capped mountains. The country is home to the Rocky Mountains, the Niagara Falls, and the Great Lakes, which offer breathtaking views and outdoor adventures for visitors.

Canada is also renowned for its diverse wildlife, including polar bears, moose, and grizzly bears, making it a popular destination for wildlife enthusiasts.

Additionally, Canada is known for its vibrant cities, such as Toronto, Vancouver, and Montreal, which offer a mix of urban amenities and natural attractions.

Overall, Canada is a country that offers a unique blend of natural beauty, cultural diversity, and urban sophistication, making it a popular destination for travelers from around the world.

Figure 10: Case study for intuitive comparisons on generic QA dataset based on Llama2-13b-base.