000 TOWARDS MITIGATING FACTUAL HALLUCINATION IN LLMs through Self-Alignment with Memory

Anonymous authors

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

Paper under double-blind review

ABSTRACT

Despite the impressive performance of Large Language Models (LLMs) across numerous tasks and widespread application in real-world scenarios, LLMs still struggle to guarantee their responses to be accurate and aligned with objective facts. This leads to **factual hallucination** of LLMs, which can be difficult to detect and mislead users lacking relevant knowledge. Post-training techniques have been employed to mitigate this issue, yet they are usually followed by a tradeoff between honesty and helpfulness, along with a lack of generalized improvements. In this paper, we propose to address it by augmenting LLM's fundamental capacity of leveraging its internal *memory*, that is, the knowledge derived from pre-training data. We introduce **FactualBench**, a comprehensive and precise factual QA dataset consisting of nearly 200k Chinese generative QA data spanning 21 domains for both evaluation and training purposes. Furthermore, we propose self-alignment with memory, i.e., fine-tuning the model via preference learning on self-generated pairwise data from FactualBench. Extensive experiments show that our method significantly enhances LLM's performance on FactualBench, with consistent improvements across various benchmarks concerning factuality, helpfulness and multiple skills. Additionally, different post-training techniques and tuning data sources are discussed to further understand their effectiveness.

1 INTRODUCTION

031 Factual hallucination occurs when large language models (LLMs) generate inaccurate or entirely 032 fabricated information in response to user queries (Zhang et al., 2023b; Huang et al., 2023). Detect-033 ing and mitigating factual hallucinations is crucial, because such errors can undermine trust in these 034 models and potentially cause significant harm to users, especially when they are used in high-stakes applications (Ji et al., 2023; Mello & Guha, 2023; Kornblith et al., 2022). However, identifying these hallucinations poses a unique challenge, as the fabricated content is often presented in a plausible 037 and convincing manner, making it difficult for users and models to recognize inaccuracies (Kaddour 038 et al., 2023; Zhang et al., 2023b). This complexity underscores the necessity of addressing factual hallucination as a critical focus for enhancing the reliability of LLMs (Tian et al., 2023). 039

040 Previous studies have explored mitigating hallucinations through post-training techniques such as 041 Supervised Fine-tuning (SFT) (Elaraby et al., 2023; Moiseev et al., 2022; Santos et al., 2022) and 042 Reinforcement Learning (RL) (Tian et al., 2023; Lin et al., 2024; Ouyang et al., 2022). However, 043 the implementation of these methods can inadvertently introduce more hallucinations (Gudibande 044 et al., 2023; Ji et al., 2023; Zhang et al., 2023b) if training on novel knowledge that model has not previously encountered (Gekhman et al., 2024; Huang et al., 2023), and create a trade-off between generating responses that are truthful and those considered helpful (Lin et al., 2024) when training 046 under inappropriate signals, for example, balancing response length and factuality in long-form tasks 047 or implicitly favoring specific response modes (Singhal et al., 2023; Sharma et al., 2024; Torabi et al., 048 2018; Kumar et al., 2024). These limitations restrict model's ability to generalize effectively to new 049 tasks and domains (Zhang et al., 2023b; Huang et al., 2023). 050

051 These challenges motivate us to enhance model factuality and facilitate generalized improvements by more effectively leveraging existing *memory* (i.e. pre-trained knowledge), a fundamental capabil-052 ity of LLMs (Zhao et al., 2023), rather than injecting new information. Specifically, we concentrate on precise closed-book question-answering (QA) task, which serves as a metric for assessing a

054 model's ability to accurately utilize its stored knowledge (Roberts et al., 2020) and has a golden cri-055 terion: the correctness of the provided answer. However, current QA datasets are often insufficient 056 and inaccurate for both comprehensive evaluation and enhancement of model performance. 057

To address the limitation, we propose FactualBench, a large-scale, multi-domain Chinese generative QA dataset designed to evaluate and improve knowledge utilization ability of LLMs. Factual-Bench is constructed from high-quality, publicly available encyclopedia entries that are commonly 060 used in pre-training corpora (Liu et al., 2024b; Ando et al., 2024), ensuring its alignment with model's existing knowledge. It comprises 181,176 benign samples across 21 domains, representing 062 a diverse set of important topics filtered by view counts. Preliminary evaluations with this dataset 063 reveal that the task is challenging for most LLMs, despite not requiring advanced skills or long-064 tailed knowledge. Interestingly, we observe that models frequently generate correct answers under high-temperature configuration (Figure 1(a)), suggesting that the knowledge is internalized but not 065 effectively utilized. This highlights a potential for improvement in factuality. 066



083 Figure 1: Left: Model's performance on FactualBench. Sampling 8 answers in high temperature 084 configuration, model is considered to have internalized the knowledge as long as one answer is correct. Orange bars indicate model's potential. Right: Model's perfomances on different benchmarks 085 before and after training. The top eight tasks are sub-dimensions of AlignBench. Green bars indicate model's improvement after self-alignment with memory. 087

The potential suggests an inappropriate utilization of memory, which could be calibrated through 089 alignment techniques. To unleash the potential of model for better factuality and generalized perfor-090 mance, we propose self-alignment with memory. Based on the train set of FactualBench that per-091 tains to memory utilization, we elicit diverse responses from the model and use these self-generated 092 outputs as labels rather than existing annotations to avoid training on their implicit answer modes. Subsequently we build pairwise data and post-train the model in preference learning (e.g., Direct 094 Preference Optimization (Rafailov et al., 2024)) to deliver a more precise bi-directional signal. 095

We evaluate model using our test set and benchmarks that assess various dimensions and abilities, 096 including CMMLU (Li et al., 2023a) for multiple-choice, HaluEval (Li et al., 2023b) for hallucination detection, TruthfulQA (Lin et al., 2022) and HalluQA (Cheng et al., 2023) for adversarial 098 robustness, and AlpacaEval (Li et al., 2023c) for helpfulness, AlignBench (Liu et al., 2023) for comprehensive abilities. As illustrated in Figue 1(b), our method achieves a unanimous and significant 100 increase in performance. We further discuss how our approach attains this effect from the relation-101 ship between memory and representation. We also conduct a series of experiments to investigate 102 how different post-training algorithms and sources of tuning data influence the training outcomes.

103 104

088

061

2 **RELATED WORKS**

105 106

Hallucination is defined as generated content that is nonsensical or unfaithful to the provided source 107 content (Ji et al., 2023; Huang et al., 2023). It presents a mutifaceted challenge Li et al. (2024b) across a wide range of tasks (Ji et al., 2023). Hallucination can be categorized into intrinsic and extrinsic parts, depending on whether the response conflicts with context itself or with the fact beyond the context (Ji et al., 2023; Huang et al., 2023), where the latter type tends to have more side effects (Zhang et al., 2023b) and is mainly focused in our work.

112 To mitigate hallucinations, several studies (Gardent et al., 2017; Wang, 2019) find that enhancing 113 the quality of pre-training data can be effective. But processing vast scale training data could be 114 time-consuming and is not be applicable to models that complete pre-training. Other approaches 115 focus on improving model's factuality through decoding strategies (Zhang et al., 2023a; Li et al., 116 2024b; Lee et al., 2022; Chuang et al., 2023), yet these strategies often increase inference complex-117 ity and have more difficulty generating fluent or diverse text (Ji et al., 2023). Additionally, some 118 methods utilize retrieval-augmented generation techniques (Nakano et al., 2021; Gou et al., 2023), but introduce significant system complexity (Tian et al., 2023) and rely on external resources. Con-119 sequently, we emphasize enhancing model factuality through post-training by directly optimizing 120 model's inherent parameters. Post-training algorithms, such as SFT and RL are frequently used to 121 improve model's instruct following ability and align model with human preferences (Xu et al., 2024; 122 Lin et al., 2024). However, SFT is considered to be sub-optimal for mitigating hallucinations due 123 to its limited generalization capabilities in out-of-distribution cases (Zhang et al., 2023b). Recent 124 RL studies improve factuality within a single task (Kang et al., 2024; Tian et al., 2023), or face 125 a trade-off between factulity and helpfulness (Lin et al., 2024). In contrast, our method achieves 126 generalized improvements across factuality, helpfulness, and comprehensive abilities. 127

The principle underlying our method is to enhance model factuality by better utilizing existing mem-128 ory. However, current datasets are often insufficient and inaccurate for evaluating and improving this 129 capability. Tasks that require long-form answers and multi-sentence reasoning (Wei et al., 2024; Min 130 et al., 2023; Joshi et al., 2017; Yang et al., 2018) are often imprecise in measuring a model's funda-131 mental knowledge utilization, as their performance is heavily influenced by instruct following and 132 complex reasoning abilities. Similarly, multiple-choice and detection tasks that rely on rule-based 133 automatic metrics (Li et al., 2023a;b; Liu et al., 2022; Mishra et al., 2024; Thorne et al., 2018) intro-134 duce significant biases into model evaluations (Lou et al., 2024). Datasets designed with adversarial 135 intent, such as TruthfulQA (Lin et al., 2022) and others (Cheng et al., 2023), effectively stimulate factual hallucinations but tend to focus on specific scenarios, thereby limiting their capacity to re-136 flect model accuracy on general, everyday questions. While datasets focused on precise generative 137 QA (Yang et al., 2015; Wang et al., 2023a; Li et al., 2024a; Yin et al., 2023; Berant et al., 2013; 138 Kwiatkowski et al., 2019) exist, they are generally constrained by small sample sizes or lack of 139 domain classification, rendering them inadequate for a comprehensive evaluation of modern LLMs. 140 These gaps in existing datasets motivate the need for more robust and large-scale resources like 141 FactualBench, which we propose to address these shortcomings. 142

143 144

145

151

152 153

3 Method

As previously discussed, we aim to mitigate model's factual hallucination while achieving general ized improvements in helpfulness and comprehensive abilities. Prior research indicates that training on existing knowledge and accurate signals is essential for optimal training outcomes. Therefore we select precise QA task, judged solely by correctness, as training task to enhance model's fundamental capability of leveraging existing memory from pre-trained corpus. The optimization goal can be mathematically formalized as follows:

$$\max_{\mathbf{x} \sim \mathcal{X}^{\text{Fact}}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [\mathcal{J}(\mathbf{x}, \mathbf{y})], \tag{1}$$

where x represents a factual question drawn from \mathcal{X}^{Fact} space. This space encompasses straightforard questions about factual information, which can be grounded to reliable sources, including but not limited to dictionaries, encyclopedias, textbooks from different domains (Wang et al., 2023b). In our study, x is selected to exclude any malicious or misleading content, and the model's response y is generated solely from question and model's internal parameters without external references. We utilize a dataset \mathbb{D}_{Test} to approximate \mathcal{X}^{Fact} . \mathcal{J} functions as an evaluation metric that outputs either 0 or 1 based on the correctness of answer y to the question x.

161 To achieve our goal, two key sub-questions should be answered: 1) How to generate a sufficient large and effective dataset to benchmark and enhance model performance on factual QA task? 2) What

methods can be employed to train the model to minimize factual hallucination, while simultaneously
 improving generalized capabilities?

165 3.1 FACTUALBENCH

167 A large-scale and comprehensive dataset containing benign and precise questions is needed for ac-168 curately evaluating and enhancing the knowledge utilization ability of models. We choose publicly available Internet encyclopedias as a reliable knowledge base as they are widely used as training corpora for LLMs (Liu et al., 2024b; Ando et al., 2024) and contain a wide range of topics across 170 various domains (Bai et al., 2024). A model-based approach is adopted to generate large amounts of 171 data costlessly and quickly. To avoid generating low-quality data (more details of such low quality 172 data can be found in Appendix A), we utilize few-shot prompts (Brown, 2020), chain-of-thought 173 (Wei et al., 2022) technologies, and apply several filtering strategies. GPT4 (Achiam et al., 2023) 174 and Baichuan model¹ are adopted in construction for their strong command-following capabilities. 175

176 177

3.1.1 CONSTRUCTION AND COMPOSITION

178 The construction pipeline for FactualBench is organized into 5 steps: 1) Entry filtering. We ini-179 tially sampled millions of entries from encyclopedia data, covering a broad spectrum of subjects and domains. To avoid testing on long-tailed knowledge, we set a view-counts threshold, and 89,658 180 encyclopedia entries remained. 2) Description filtering. Model's performance tends to decline as 181 context length increases (Liu et al., 2024a; Sun et al., 2023; Li et al., 2024c). Excessively lengthy 182 description may lead to low quality responses. Conversely, a overly brief description lacks suffi-183 cient factual information. To balence this, we filter out descriptions shorter than 100 characters and 184 truncated those exceeding 800 characters. 64315 entries remained after this process. 3) Question 185 generation. We instruct GPT4 to generate at most 3 precise questions based on each truncated description. For each question, GPT4 is also required to provide 1 standard answer and 3 misleading 187 incorrect answers. To ensure adherence to our instructions, we included two examples in the prompt 188 as few shots. Totally 192,927 QA data are generated. 4) Question classification. Following ques-189 tion generation, a domain classifier, fine-tuned on Baichuan model, is employed to categorize all questions into 20 distinct domains. Questions that don't fit into any domain are uniformly classi-190 fied as others. 5) Question filtering. We query GPT4 once more to filter out low-quality questions. 191 Each question is assessed without corresponding encyclopedia description and GPT4 is instructed to 192 judge whether question is low-quality through a step-by-step reasoning. Finally 181,176 questions 193 are reserved, among which high-quality data accounts for 90% under human assessment. The com-194 plete prompts utilized throughout the generation pipeline can be found in Appendix B.1.1. Table 3 195 presents a sample from our dataset, while additional examples are provided in Appendix B.2. 196

Table 1: Samples distribution of Factualbench.

Table 2: Models performances on Factual-Bench rated by GPT4.

						Denen faicu by OF 14.		
199	Domain	中文名	Test	Training	Total	5		
200	film&entertainment	影视娱乐	201	54489	54690	Model	Proficient lang	100
	eduaction&training	教育培疗	161	3703	3864	WIOdel	r toncient lang.	Acc.
201	physics, chemistry, mathematics&biology	数理化生	201	9189	9390	Daishuan 1	CN	40.24
	history&traditional culture	历史国学	202	18108	18310	Balchuani	CN	48.24
202	biography	人物百科	201	11844	12045	Baichuan2	CN	55.37
000	politics&law	政治法律	175	6368	6453	Owen1.5-7B	CN	48.87
203	economics&management	经济管理	160	4543	4703	Owen2 7P	CN	56 27
004	computer science	计算机科学	201	6253	6454	Qwell2-7B	CN	50.27
204	medical	医学	167	7073	7240	Llama-3-8B	EN	39.11
205	sociology&humanity	社会人文	199	8503	8702		<i>a</i> .	(7.50
200	agriculture, forestry, fisheries&allied industries	农林牧渔	153	3728	3881	Baichuan3	CN	67.50
206	astronomy&geography	天文地理	160	3896	4056	Yi-34B	CN	67.30
200	sports&tourism	运动旅游	157	4869	5026	Command-R	FN	54 30
207	digital&automotive	数码汽车	176	3887	4063	L1 2 70D	EN	10.65
201	industrial engineering	工业工程	172	3283	3455	Llama-3-70B	EN	49.65
208	military&war	军武战争	151	2569	2720	Qwen2-72B	CN	73.71
200	slang&memes	网词网梗	151	529	680			
209	work&life	工作生活	174	5853	6027	Baichuan4	CN	75.07
100	high technology	高新科技	150	310	460	Command-R+ 104B	EN	60.17
210	religion&culture	信仰文化	150	510	660	DeenSeelt v2 0627 MeE 226P	CN	75.62
	others	其他	/	18207	18207	DeepSeek-v2-002/ MOE-250B	CN	75.02
211	total	/	3462	177714	181176	GP14-0125-preview	EN	65.71

211 212 213

214

215

197

To establish an efficient benchmark, we randomly select 3,462 samples as the test set, while the remaining 177,714 samples comprised the training set. This selection is based on encyclopedia

¹https://www.baichuan-ai.com/

entries instead of questions, ensuring all data in test set seperate from training set. Each domain has a similar number of questions in test set and entries containing *others* domain questions are excluded from selection. We manually rephrase unclear questions to maintain the quality of test set.
The distribution of FactualBench is presented in Table 1.

Table 3: Each sample of FactualBench contains a question, a corresponding standard answer, three misleading incorrect answers and domain it belongs to. English translation is for reference only.

Question	第一台微波量子放大器是在哪一年制成的?	In which year was the first microwave quantum amplifier made?
Standard Answer	第一台微波量子放大器是在1954年制成的。	The first microwave quantum amplifier was made in 1954.
Wrong Answer1 Wrong Answer2 Wrong Answer3	第一台微波量子放大器是在1958年制成的。 第一台微波量子放大器是在1960年制成的。 第一台微波量子放大器是在1962年制成的。	The first microwave quantum amplifier was made in 1958. The first microwave quantum amplifier was made in 1960. The first microwave quantum amplifier was made in 1962.
Domain	高新科技	high technology

232 3.1.2 EVALUATION

220

221

222

Similar to previous works (Liu et al., 2023; Zheng et al., 2024), a robust model-based approach is employed to expedite the assessment process. Given that rule-based automatic metrics such as ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002) exhibit significant biases in evaluation (Lou et al., 2024), we assess the correctness of answer at semantic-level. The question, model answer, and standard answer are provided for judgement and evaluator is supposed to focus solely on the content that directly addressed the question and determine whether the model's answer aligns with the standard answer.

240 In our present work, the evaluator judges an answer as correct (i.e. $\mathcal{J}(\mathbf{x}, \mathbf{y}) = 1$ in equation 1) only 241 when model answers the question and the answer matches the standard answer. It is reasonable since 242 model is expected to have been trained on relevant data and therefore should possess the knowledge. 243 Moreover, this knowledge is not long-tailed, but rather high frequency viewed. Additionally, we an-244 alyze the responses from different tested models and find that the proportion of evasive answers (e.g. 245 "I don't know") is only approximately 1%. To enhance judgment accuracy, we instruct the evaluator 246 to provide analysis and include several examples in the prompt. GPT4-0125-preview is chosen as 247 the evaluator, achieving a 96% consistency rate with human, thus validating the effectiveness and accuracy of our evaluation process. The prompt used for evaluation is detailed in Appendix B.1.2. 248

We evaluate 14 popular open and closed source LLMs on FactualBench, including Baichuan series (Yang et al., 2023), Qwen series (Yang et al., 2024; Bai et al., 2023), Llama-3 series (AI@Meta, 2024), Yi (AI et al., 2024), Command-R series (Gomez, 2024a;b), DeepSeek (DeepSeek-AI, 2024), and GPT4 (Achiam et al., 2023). We prioritize the chat/instruct versions of these models. Detailed settings are shown in Appendix C.1 and the results are listed in Table 2. The accuracy of LLMs on our benchmark ranges from 39.11% to 75.62%, indicating that most models still have lacks even in basic factual QA task. Detailed results for each domain and analysis are shown in Appendix D.

257 3.1.3 POTENTIAL IN MODEL

258 For questions where model response incorrectly, we observe that it can still generate correct answers 259 when allowed greater diversity in its outputs. Taking Baichuan1 as an example, we encourage the 260 model to generate varied answers by increasing the generation temperature. By sampling model's 261 responses 8 times under this condition (which we refer as *high temperature BO8*), as opposed to 262 the standard inference condition (termed low temperature BO1). We consider the model to possess 263 relevant information and the ability to utilize it if at least one of the generated answers is correct. As 264 illustrated in Figure 1(a), comparing the accuracy of BO8 and BO1, we find a significant portion of the model's capability have not been fully stimulated, i.e., potential. 265

266 267

268

- 3.2 Self-alignment with memory
- Transformer (Vaswani et al., 2017) achitecture LLMs are trained to solve next-word-prediction tasks follows a statistical paradigm (Arora & Goyal, 2023). As probabilistic models, LLMs can generate



Table 4: The construction of tuning set.

Questions	Correctly answered questions	Correct answers
24,000 177,714	15,489 115,798 (SFT1)	/ 489,357 (SFT2)
Questions	Valid questions	DPO pairs

Figure 2: We align model with self-generated data on task related to memory utilization.

275

276

277 278

293

298 299

306

307

314 315

316

diverse answers to the same question based on sampling from a distribution, derived from extensive
training data. The model's ability to provide correct answer at high temperature indicates a extent
of internalization of relevant knowledge. However, an incorrect distribution will result in a high
probability of sampling incorrect answers, which can be caused by insufficient or wrong alignment.
Since knowledge utilization is a fundamental ability of LLMs (Zhao et al., 2023), LLMs could have
a generalized improvement if better alignment is achieved.

Ξ

285 To stimulate potential and enhance knowledge utilization ability, we introduce self-alignment with memory. According to Figure 2, we first collect model's high temperature BO8 answers on the 286 training set of FactualBench, a task exclusively evaluate LLM's memory on pre-training knowledge. 287 Then we evaluate answers' correctness, but only a weaker evaluator is needed instead of GPT4. 288 We choose model's self-generated data as tuning labels to prevent the risk of model hallucination 289 exacerbated by fine-tuning on new knowledge (Gekhman et al., 2024) or learning implicit response 290 patterns from external annotations (Kumar et al., 2024). Among BO8 answers, we construct a maximum of 8 pairwise data (x, y_w, y_l) per question, using correct answers as chosen labels and 291 the wrong ones as rejected labels, which can be formulated by the following constraint conditions: 292

$$\mathbf{x} \sim \mathbb{D}^{\text{train}}; \ \mathbf{y} \sim \pi_{\text{ref}}(\cdot | \mathbf{x}); \ \mathcal{J}(\mathbf{x}, \mathbf{y}_w) = 1; \ \mathcal{J}(\mathbf{x}, \mathbf{y}_l) = 0.$$
(2)

In this way, we can quickly generate tuning set $\mathbb{D}^{\text{tuning}}$ containing massive data without human involvement. Then fine-tune the model on this tuning data in preference way, e.g. Direct Preference Optimization (DPO) (Rafailov et al., 2024) for a precise control through bi-directional signals. The DPO loss is defined as follows:

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_{w}, \mathbf{y}_{l}) \sim \mathbb{D}^{\text{tuning}}} \left[\log \sigma(\beta \log \frac{\pi_{\theta}(\mathbf{y}_{w} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{w} | \mathbf{x})} - \beta \log \frac{\pi_{\theta}(\mathbf{y}_{l} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{l} | \mathbf{x})}) \right],$$
(3)

where π_{θ} is the optimal policy, π_{ref} is the policy before optimization. σ denotes for sigmoid function, and β is a hyperparameter.

The construction of the tuning data is illustrated in Table 4. We randomly select a subset of 24k questions for an early checkpoint and a question is considered to be valid only if it receives both correct and incorrect answers.

4 EXPERIMENTS AND ANALYSIS

In this section, we present the training result obtained using our method. Comparing the effectiveness of SFT and DPO, and tuning data from different sources, we observe that: 1) Training on precise factual QA task that related to pre-training knowledge can lead to generalized improvements on other tasks; 2) Training on pairwise data has an advantage over pointwise data in the task; 3) Training effect is enhanced when tuning on self-generated data. For pairwise data, it is crucial that labels from both directions adhere to the same distribution to avoid being hacked easily.

4.1 Setup

We use Baichuan1 as the experimental model. Beside FactualBench, we also choose 6 other benchmarks to evaluate generalized improvements on other factual tasks, helpfulness and general capabilities: TruthfulQA (Lin et al., 2022) for English, HalluQA (Cheng et al., 2023) for Chinese, CMMLU (Li et al., 2023a) for multiple-choice task, HaluEval (Li et al., 2023b) for detection task, AlpacaEval (Li et al., 2023c) for helpfulness, and AlignBench (Liu et al., 2023) for comprehensive abilities. We introduce the details of benchmarking in Appendix C.2.

323 Traditional SFT (SFT0) and DPO (DPO0) are conducted as baselines, whose labels are directly provided by GPT4 annotations in FactualBench dataset. Besides, we use self-generated data to train

Table 5: Training results after SFT and DPO tuning on Baichuan1. ¹AlpacaEval is calculated on model answers' win rate against the model before training. ²When calculating Avg Δ , we multi-ply the score of AlignBench by 10 to align the accuracy of other benchmarks. AlpacaEval shows a relative win rate hence is not concluded in Avg Δ . Red underline indicates an improvement af-ter training; Bold font indicates the best performance on the metric; Red font indicates an overall improvement.

Baichuan1	cho	sen/sft label	rejected lat	el Factu	alBench(3642)/acc.	TruthfulQ	A/acc.	HalluQA/acc.	CMMLU/acc.	HaluEval/a	cc. Align	Bench/score	AlpacaEval/w	in rate 1	Avg Δ^2
Chat		/	/		48.24	30.2	3	32.00	48.85	50.35	-	5.03	50.00		/
SFT0		dataset	/		55.86(+7.62)	21.30(-	3.93)	22.44(-9.56)	49.58(+0.73)	12.40(-37.	95) 3.1	73(-1.30)	26.65		-10.18
SFT1 SFT2		BO8 BO8	/		51.33(+3.09) 52.37(+4.13)	31.46(+ 28.76(-	1.23) .47)	30.00(-2.00) 26.44(-5.56)	48.78(-0.07) 50.15(+1.30)	55.73(+5.) 53.90(+3.)	18) <u>5.0</u> (5) <u>5.0</u>	04(+0.01) 03(-0.00)	37.58 31.06		+1.29 +0.32
DPO0		dataset	dataset		49.08(+0.84)	28.89(-	.34)	19.78(-12.22)	50.70(+1.85)	54.89(+4.3	i <u>4)</u> 4.8	82(-0.21)	39.07		-1.40
DPO1 (ours) DPO2 (ours)		BO8 BO8	BO8 BO8		57.37(+9.13) 58.29(+10.05)	33.78(+ 35.86(+	3.55) 5.63)	38.44(+6.44) 38.89(+6.89)	50.13(+1.28) 50.92(+2.07)	50.63(+0.1 52.05(+1.1	28) 5.3 70) 5.3	80(+0.27) 88(+0.35)	54.84 63.99		+3.90 + 4.97
SFT then DPO SFT and DPO		BO8 BO8	BO8 BO8		54.74(+6.50) 57.16(+8.92)	37.33(+ 34.76(+	7.10) 4.53)	36.67(+4.67) 38.22(+6.22)	50.72(+1.87) 50.78(+1.93)	54.02(+3.0 52.31(+1.9	5.0 5.1 5.1	07(+0.04) 13(+0.10)	54.53 63.91		+4.03 +4.09
Baichuan1	Avg.	Professional K	nowledge	Mathematic	s Fundamental Langu	age Ability	Logica	Reasoning Ad-	anced Chinese Uno	lerstanding V	Vriting Abilit	y Task-orier	nted Role Play 0)pen-endeo	Questions
Chat	5.03	5.34		2.71	5.57			3.20	5.86		6.32		6.33	6.	53
SFT0	3.73	4.48(-0.	86)	2.62(-0.09)	4.79(-0.78	i)	2.7	5(-0.45)	5.08(-0.78)		3.24(-3.08)	3.77	7(-2.56)	3.76(-	2.87)
SFT1 SFT2	5.04 5.03	5.78(+0. 5.46(+0.	.44) .12)	2.59(-0.12) 2.88(+0.17)	5.47(-0.10 5.60(+0.03	1) 3)	3.3	0(+0.10) 5(+0.05)	5.66(-0.20) 5.57(-0.29)		6.11(-0.21) 6.19(-0.13)	6.25 6.17	5(-0.08) 7(-0.16)	6.58(- 6.63(0.05) 0.00)
DPO0	4.82	4.67(-0.	67)	2.60(-0.11)	5.53(-0.04	-)	3.30	0(+0.10)	5.50(-0.36)		6.40(+0.08)	6.17	7(-0.16)	6.00(-	0.63)
DPO1 (ours) DPO2 (ours)	5.30 5.38	5.92(+0. 6.25(+0.	.58) . 91)	3.02(+0.31) 3.03(+0.32)	5.66(+0.09 5.76(+0.19))))	3.3 3.5	7(+0.17) 5(+0.35)	5.97(+0.11) 6.12(+0.26)		6.53(+0.21) 6.52(+0.20)	6.55 6.36	5(+0.22) 5(+0.03)	6.79(+ 6.79(+	-0.16) -0.16)
SFT then DPO SFT and DPO	5.07 5.13	5.57(+0. 5.60(+0.	.23) .26)	2.66(-0.05) 2.79(+0.08)	5.53(-0.04 5.57(+0.00	-)))	3.0 3.1	1(-0.19) 6(-0.04)	6.00(+0.14) 6.05(+0.19)		$\frac{6.33(+0.01)}{6.17(-0.15)}$	6.32 6.41	2(-0.01) (+0.08)	6.92(+ 7.16(+	-0.29) -0.53)

SFT in two settings: one label for each question (SFT1) and all correct answers as labels for each question (SFT2); And we train DPO using our self-alignment with memory in different data sizes (DPO1 and DPO2). Some works claim that fusing DPO loss with SFT loss can help mitigate the overoptimization on rejected labels (He et al., 2024; Liu et al., 2024c), we donate this process as SFT and DPO; And additional SFT training before DPO on the preference tuning set can reduce distribution shift issue therefore help training (Xu et al., 2024), which we donate as SFT then DPO. More training details are shown in Appendix C.3.

4.2 **RESULTS**

We show the results in Table 5. Our method (DPO2) achieves unanimous improvement on all benchmarks and sub-dimensions of AlignBench, with an average improvement of 4.97% on 6 benchmarks, and a helpful win rate of 63% against Chat model, demonstrating the effectiveness of our method. According to Figure 3, our method stimulates the potential of model, and the improvement is mainly sourced from existing memory rather than new knowledge.



Figure 3: Performance on FactualBench after self-alignment with memory. Left: The accuracy of low temperature BO1, SFT, DPO (ours), and high temperature BO8 on different domains (each domain is represented by its first 5 letters). Right: Venn graph showing the correct answer distribution between low temperature BO1, DPO (ours), and high temperature BO8.

In contrast, traditional training methods that directly use the ground truth annotations exhibit fluc-tuations in performance on new tasks: hallucinate more on TruthfulQA and HalluQA, and have decline in fundamental abilities and helpfulness. SFT0 even shows a serious decline in instruction following on HaluEval. Since annotations in dataset are relatively short and concise, they could
work as biased signals, which emphasizing the advantage of self-generated data. Comparing SFT
and DPO, we find that pairwise data will lead to a greater improvement on training performance
(although DPO1 has less tuning data than SFT1), indicating that using unidirectional signal brought
by SFT is still insufficient for our task.

We use different sizes of data for DPO training and present the training results in Figure 4. As the size of training data increases, the overall improvement, measured by Avg Δ , of the model shows an approximate logarithmic curve, indicating early training on our training set can already effectively improve the model's ability.

In addition, comparing SFT1 and SFT2, we observe that using more correct labels on the same question in SFT doesn't improve the training effect, which limits the effective tuning data size of SFT when the number of training questions is fixed. We also notice that neither fusing SFT loss in DPO loss nor conduct a SFT training before DPO training improves the training results. Since the pairwise data are all sampled from the model itself, there will be little distribution shift and is difficult to overoptimize solely on rejected lables, which also indicates the training robustness of the data constructed by our method.

So far, it still need to be answered why self-alignment with memory effective for mitigating factual hallucination and could lead to generalized improvement.



Figure 4: Fitted logarithmic scaling law for training data size of DPO training on Bachuan1.
The improvement shows a decreasing marginal benefit trend with the increase of training data size.



Figure 5: Loss comparison between different sources of tuning data.

4.3 ANALYSIS

394

395

417 418

419

420 In transformer architecture LLMs, Attention layers and Multilayer Perceptron (MLP) layers extract 421 useful features from input (Jiang et al., 2024), in which MLP layers are regarded to implement a 422 lookup table for factual recall to output multi-token embedding with related information (Nanda et al., 2023) and provide factual knowledge (Dai et al., 2022). Inspired by the skill graph (Arora & 423 Goyal, 2023), we summarize the role of MLP as a memory graph, with edge connecting different 424 input contexts and internalized knowledge. A more precise memory graph leads to better knowledge 425 utilization, and therefore results in accurate answer distribution. The task we choose is highly related 426 to knowledge utilization, providing a direct optimization signal for the ability, a fundamental ability 427 that could affect diverse tasks. 428

In contrast to the single signal provided in SFT, preference learning like DPO provides a bidirectional signal, which can weaken useless or wrong edge as well as connect or strengthen an needed edge, a more precise control than SFT. Besides, there could be diverse data samples on the same input, making training process more robust. Compared to external label, a self-generate label

has the same answer mode as model itself, avoiding overoptimizing on new answer pattern. And we could also explain why a model fine-tuned on new knowledge it doesn't possess will cause more hal-lucinations: Since there is no corresponding knowledge in memory graph, model cannot learn how to obtain the correct distribution but an overfitting on the label through building incorrect edges.

A better knowledge utilization leads to better representation. Recent work (Huh et al., 2024) shows that representation alignment increases with LLM's scale and performance. We choose Qwen2-72B-Instruct (Bai et al., 2023) as a good representation function, and the last hidden state value as representation. Our experiments prove that after DPO training, Baichuan1 has further aligned with Qwen2-72B, indicating a better representation ability is achieved. We use mutual nearest-neighbor metrics (Huh et al., 2024) (we explain the calculation process of the metric in the Appendix E) to evaluate alignment between two models, set k = 10, and randomly select 200 data points to calculate.

We show the alignment change before and after training based on FactualBench in Figure 6, and alignment changes based on another four benchmarks in Figure 7. Our method achieves higher align-ment with Qwen2-72B on four benchmarks and most domains in FactualBench, indicates model has gained stronger representation ability, and result in a better response accuracy. Since different tasks have different difficulties and features, the alignment cannot accurately predict the performance of the model. However, it still works as a signal, reflecting the trend of the model performance changing. Comparing with SFT, DPO has higher alignment degree on average, proving that bi-directional signal of DPO is better for model to achieve generalized improvements.



Figure 6: The change of Baichuan1's alignment with Qwen2-72B-Instruct on FactualBench after training. Each point represents a domain. Left: DPO (ours); Right: SFT.



Figure 7: The changes of Baichuan1's alignment with Qwen2-72B-Instruct on four benchmarks after training. From left to right: TruthfulQA, HalluQA, HaluEval, AlignBench. Acc. and score are calculated on selected samples.

4.4 ABLATION STUDIES

We conduct several ablation experiments to investigate the impact of different sources of tuning data. Additionally, we experiment with Qwen2-7B-Instruct to validate the effectiveness of our method, but only use the subset of 24,000 questions as early training can already effectively improve model's Table 6: Training results of ablation studies on Baichuan1 and Qwen2-7B. The default BO8 data are sampling from the model to be trained. BC.BO8 denotes for BO8 sampling from Baichuan1. w/ desc denotes for BO1 with description.

Baichuan l	chosen/sft label	rejected	FactualBench/acc.	TruthfulQA/acc.	HalluQA/acc.	CMMLU/acc.	HaluEval/acc.	AlignBench/score	AlpacaEval/win rate	Avg Δ
Chat	/	/	48.24	30.23	32.00	48.85	50.35	5.03	50.00	/
SFT1	BO8	/	51.33(+3.09)	31.46(+1.23)	30.00(-2.00)	48.78(-0.07)	55.73(+5.38)	5.04(+0.01)	37.58	+1.29
SFT3	w/ desc	/	55.63(+7.39)	36.60(+6.37)	27.11(-4.89)	51.39(+2.54)	10.40(-39.95)	4.47(-0.56)	36.96	-5.69
SFT0	dataset	/	55.86(+7.62)	21.30(-8.93)	22.44(-9.56)	49.58(+0.73)	12.40(-37.95)	3.73(-1.30)	26.65	-10.18
DPO2 (ours)	BO8	BO8	58.29(+10.05)	35.86(+5.63)	38.89(+6.89)	50.92(+2.07)	52.05(+1.70)	5.38(+0.35)	63.99	+4.97
DPO3	w/ desc	BO8	18.17(-30.07)	13.10(-17.13)	9.33(-22.67)	48.05(-0.80)	48.57(-1.78)	4.07(-0.96)	32.80	-13.67
DPO4	dataset	BO8	5.40(-42.84)	3.92(-26.31)	1.56(-30.44)	46.85(-2.00)	40.10(-10.25)	3.28(-1.75)	19.07	-21.56
DPO0	dataset	dataset	49.08(+0.84)	28.89(-1.34)	19.78(-12.22)	50,70(+1.85)	54.89(+4.54)	4.82(-0.21)	39.07	-1.40
Qwen2-7B	chosen/sft label	rejected	FactualBench/acc.	TruthfulQA/acc.	HalluQA/acc.	CMMLU/acc.	HaluEval/acc.	AlignBench/score	AlpacaEval/win rate	Avg Δ
Qwen2-7B Instruct	chosen/sft label /	rejected /	FactualBench/acc. 56.27	TruthfulQA/acc. 52.75	HalluQA/acc. 46.44	CMMLU/acc. 80.85	HaluEval/acc. 52.30	AlignBench/score 6.69	AlpacaEval/win rate 50.00	Avg Δ
Qwen2-7B Instruct SFT4	chosen/sft label / BO8	rejected / /	FactualBench/acc. 56.27 55.43(-0.84)	TruthfulQA/acc. 52.75 50.31(-2.44)	HalluQA/acc. 46.44 45.56(-0.88)	CMMLU/acc. 80.85 80.22(-0.63)	HaluEval/acc. 52.30 53.70(+1.40)	AlignBench/score 6.69 6.63(-0.06)	AlpacaEval/win rate 50.00 44.22	Avg Δ / -0.66
Qwen2-7B Instruct SFT4 SFT5	chosen/sft label / BO8 BC.BO8	rejected / / / /	FactualBench/acc. 56.27 55.43(-0.84) 49.97(-6.30)	TruthfulQA/acc. 52.75 50.31(-2.44) 29.87(-22.88)	HalluQA/acc. 46.44 45.56(-0.88) 24.67(-21.77)	CMMLU/acc. 80.85 80.22(-0.63) 77.49(-3.36)	HaluEval/acc. 52.30 53.70(+1.40) 42.05(-10.25)	AlignBench/score 6.69 6.63(-0.06) 4.97(-1.72)	AlpacaEval/win rate 50.00 44.22 15.03	Avg Δ / -0.66 -13.63
Qwen2-7B Instruct SFT4 SFT5 SFT6	chosen/sft label / BO8 BC.BO8 dataset	rejected / / / /	FactualBench/acc. 56.27 55.43(-0.84) 49.97(-6.30) 50.38(-5.89)	TruthfulQA/acc. 52.75 50.31(-2.44) 29.87(-22.88) 19.58(-33.17)	HalluQA/acc. 46.44 45.56(-0.88) 24.67(-21.77) 21.11(-25.33)	CMMLU/acc. 80.85 80.22(-0.63) 77.49(-3.36) 79.85(-1.00)	HaluEval/acc. 52.30 53.70(+1.40) 42.05(-10.25) 9.69(-42.61)	AlignBench/score 6.69 6.63(-0.06) 4.97(-1.72) 3.56(-3.13)	AlpacaEval/win rate 50.00 44.22 15.03 7.20	Avg Δ / -0.66 -13.63 -23.22
Qwen2-7B Instruct SFT4 SFT5 SFT6 DPO5 (ours)	chosen/sft label / BO8 BC.BO8 dataset BO8	rejected / / / / BO8	FactualBench/acc. 56.27 55.43(-0.84) 49.97(-6.30) 50.38(-5.89) 58.81(+2.54)	TruthfulQA/acc. 52.75 50.31(-2.44) 29.87(-22.88) 19.58(-33.17) 54.47(+1.72)	HalluQA/acc. 46.44 45.56(-0.88) 24.67(-21.77) 21.11(-25.33) 49.78(+3.34)	CMMLU/acc. 80.85 80.22(-0.63) 77.49(-3.36) 79.85(-1.00) 82.15(+1.30)	HaluEval/acc. 52.30 53.70(+1.40) 42.05(-10.25) 9.69(-42.61) 54.00(+1.70)	AlignBench/score 6.69 6.63(-0.06) 4.97(-1.72) 3.56(-3.13) 6.96(+0.27)	AlpacaEval/win rate 50.00 44.22 15.03 7.20 58.26	Avg Δ / -0.66 -13.63 -23.22 +2.22
Qwen2-7B Instruct SFT4 SFT5 SFT6 DPO5 (ours) DPO6	chosen/sft label / BO8 BC.BO8 dataset BO8 BC.BO8	rejected / / / BO8 BC.BO8	FactualBench/acc. 56.27 55.43(-0.84) 49.97(-6.30) 50.38(-5.89) 58.81(+2.54) 58.17(+1.90)	TruthfulQA/acc. 52.75 50.31(-2.44) 29.87(-22.88) 19.58(-33.17) 54.47(+1.72) 53.86(+1.11)	HalluQA/acc. 46.44 45.56(-0.88) 24.67(-21.77) 21.11(-25.33) 49.78(+3.34) 46.67(+0.23)	CMMLU/acc. 80.85 80.22(-0.63) 77.49(-3.36) 79.85(-1.00) 82.15(+1.30) 80.14(-0.71)	HaluEval/acc. 52.30 53.70(+1.40) 42.05(-10.25) 9.69(-42.61) 54.00(+1.70) 52.26(-0.04)	AlignBench/score 6.69 6.63(-0.06) 4.97(-1.72) 3.56(-3.13) 6.96(+0.27) 6.71(+0.02)	AlpacaEval/win rate 50.00 44.22 15.03 7.20 58.26 39.19	Avg Δ / -0.66 -13.63 -23.22 +2.22 +0.45
Qwen2-7B Instruct SFT4 SFT5 SFT6 DPO5 (ours) DPO6 DPO7	chosen/sft label / BO8 BC.BO8 dataset BO8 BC.BO8 dataset	rejected / / / BO8 BC.BO8 dataset	FactualBench/acc. 56.27 55.43(-0.84) 49.97(-6.30) 50.38(-5.89) 58.81(+2.54) 58.17(+1.90) 55.75(-0.52)	TruthfulQA/acc. 52.75 50.31(-2.44) 29.87(-22.88) 19.58(-33.17) 54.47(+1.72) 53.86(+1.11) 52.14(-0.61)	HalluQA/acc. 46.44 45.56(-0.88) 24.67(-21.77) 21.11(-25.33) 49.78(+3.34) 46.67(+0.23) 46.22(-0.22)	CMMLU/acc. 80.85 80.22(-0.63) 77.49(-3.36) 79.85(-1.00) 82.15(+1.30) 80.14(-0.71) 80.77(-0.08)	HaluEval/acc. 52.30 53.70(+1.40) 42.05(-10.25) 9.69(-42.61) 54.00(+1.70) 52.26(-0.04) 51.70(-0.60)	AlignBench/score 6.69 6.63(-0.06) 4.97(-1.72) 3.56(-3.13) 6.96(+0.27) 6.71(+0.02) 6.50(-0.19)	AlpacaEval/win rate 50.00 44.22 15.03 7.20 58.26 39.19 36.06	Avg Δ / -0.66 -13.63 -23.22 +2.22 +0.45 -0.65

abilities. In this section, we introduce another data source: *BO1 with description*, which provides reference description to assist in questioning the model under a low temperature configuration. Since the standard answer is contained in description, BO1 with description answers are basically correct, and we use them as chosen/sft label for DPO and SFT training. Despite being generated by the same model, the distribution of BO1 with description answer still differs significantly from the model's own distribution due to the varying input contexts. The settings and the training results are shown in Table 6.

For both SFT and DPO, self-generated data yield relative better results. Comparing SFT1, SFT3, and SFT0, although SFT3 and SFT0 exhibit greater improvements on FactualBench, they have significant declines in instruct following (HaluEval) and comprehensive ability (Alignbench). Comparing SFT4, SFT5, and SFT6, the Avg Δ of SFTs all decrease, suggesting that achieving generalized improvements and stimulating a model's potential via SFT is challenging. As for DPO, DPO2 and DPO5 utilizing self-generated data both achieve the best performances and consistent improvements across all benchmarks, underscoring the effectiveness of bi-directional signals.

516 For DPO, training data from alternative sources could still yield positive effects. But chosen 517 and rejected labels should be sampled from the same distribution. Comparing DPO5 and DPO6, 518 which are trained on Qwen2 with data generated by Qwen2 and Baichuan1 both achieve improve-519 ments on FactualBench and Avg Δ . However, when chosen and rejected labels are sourced from 520 different distributions (DPO3, DPO4), the training process is quickly hacked, resulting in undesir-521 able parameter changes. According to Figure 5, the loss curve shows a straightly downward trend at the begining, leading to a deterioration of basic conversational abilities, characterized by repetitive 522 523 and incoherent responses. In cases where chosen and rejected labels are all from dataset, although they are not strictly from the same distribution, it still difficult to be hacked, and therefore there are 524 no significant performance degradation after training. 525

Detailed training results on FactualBench and AlignBench can be found in Appendix F.

526 527 528

529

504

505

506

507

508

5 CONCLUSION

530 In this article, we aim to mitigate factual hallucination and achieve generalized improvement in 531 model performance. We select factual QA as our training task to enhance model's ability to utilize 532 its memory, i.e. existing knowledge derived from pre-training data. We first extract knowledge from 533 encyclopedia to construct a large-scale, multi-domain Chinese factual QA dataset FactualBench. 534 Based on FactualBench, we observe that model still possesses significant potential for knowledge utilization. Consequently we propose self-alignment with memory: construct self-generated tuning 536 data on FactualBench and train the model using DPO loss. Our method significantly improves 537 model's performance on our benchmark as well as multiple other open-source benchmarks that evaluate factuality, helpfulness, and comprehensive skills. We attribute the effectiveness of our 538 method is originated from better representation ability. Finally, we establish the necessity of bidirectional signals and self-generated data through a series of ablation experiments.

540 REFERENCES

551

552

553

558

559 560

568

577

578

579

580 581

582

583

584

585

586

587 588

589

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024. URL https://arxiv.org/abs/2403.04652.
 - AI@Meta. Llama 3 model card. https://github.com, 2024. URL https://github.com/ meta-llama/llama3/blob/main/MODEL_CARD.md.
- Kenichiro Ando, Satoshi Sekine, and Mamoru Komachi. Wikisqe: A large-scale dataset for sentence quality estimation in wikipedia. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17656–17663, 2024.
 - Sanjeev Arora and Anirudh Goyal. A theory for emergence of complex skills in language models. arXiv preprint arXiv:2307.15936, 2023.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Yuelin Bai, Xinrun Du, Yiming Liang, Yonggang Jin, Ziqiang Liu, Junting Zhou, Tianyu Zheng,
 Xincheng Zhang, Nuo Ma, Zekun Wang, et al. Coig-cqia: Quality is all you need for chinese instruction fine-tuning. *arXiv preprint arXiv:2403.18058*, 2024.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1533–1544, 2013.
- 576 Tom B Brown. Language models are few-shot learners. arXiv preprint ArXiv:2005.14165, 2020.
 - Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin Wang, Xiangyang Liu, Mozhi Zhang, Junliang He, Mianqiu Huang, Zhangyue Yin, Kai Chen, et al. Evaluating hallucinations in chinese large language models. *arXiv preprint arXiv:2310.03368*, 2023.
 - Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
 - Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 8493–8502, 2022.
 - DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos
 Guestrin, Percy S Liang, and Tatsunori B Hashimoto. Alpacafarm: A simulation framework for
 methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

- Mohamed Elaraby, Mengyin Lu, Jacob Dunn, Xueying Zhang, Yu Wang, Shizhu Liu, Pingchuan Tian, Yuping Wang, and Yuxuan Wang. Halo: Estimation and reduction of hallucinations in open-source weak large language models. *arXiv preprint arXiv:2308.11764*, 2023.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. Creating training corpora for nlg micro-planning. In *55th Annual Meeting of the Association for Computational Linguistics*, *ACL 2017*, pp. 179–188. Association for Computational Linguistics (ACL), 2017.
- Zorik Gekhman, Gal Yona, Roee Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan
 Herzig. Does fine-tuning llms on new knowledge encourage hallucinations? *arXiv preprint arXiv:2405.05904*, 2024.
- Aidan Gomez. Command r: Retrieval-augmented generation at production scale. https://cohere.com, 2024a. URL https://cohere.com/blog/command-r.
- Aidan Gomez. Introducing command r+: A scalable llm built for business. https://cohere.com,
 2024b. URL https://cohere.com/blog/command-r-plus-microsoft-azure.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen.
 Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2023.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*, 2023.
- Qianyu He, Jie Zeng, Qianxi He, Jiaqing Liang, and Yanghua Xiao. From complex to simple:
 Enhancing multi-constraint complex instruction following ability of large language models. *arXiv* preprint arXiv:2404.15846, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023.
- Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.

632

633

- Yiding Jiang, Christina Baek, and J Zico Kolter. On the joint interaction of models, data, and features. In *The Twelfth International Conference on Learning Representations*, 2024.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly
 supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meet- ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611,
 2017.
- Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and
 Robert McHardy. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*, 2023.
- Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. Unfamiliar finetuning examples control how language models hallucinate. *arXiv preprint arXiv:2403.05612*, 2024.
- Aaron E Kornblith, Chandan Singh, Gabriel Devlin, Newton Addo, Christian J Streck, James F
 Holmes, Nathan Kuppermann, Jacqueline Grupp-Phelan, Jeffrey Fineman, Atul J Butte, et al.
 Predictability and stability testing to assess clinical decision instrument performance for children after blunt torso trauma. *PLOS digital health*, 1(8):e0000076, 2022.

648 649 650	Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, et al. Training language models to self-correct via reinforcement learning. <i>arXiv preprint arXiv:2409.12917</i> , 2024.
651	
652	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
653	Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a
654	benchmark for question answering research. Transactions of the Association for Computational
655	Linguistics, 7:453–466, 2019.
656	Naveon Lee Wei Ping Peng Xu Mostofa Patwary Pascale N Fung Mohammad Shoeybi and Bryan
657	Catanzaro. Factuality enhanced language models for open-ended text generation. Advances in
658 659	Neural Information Processing Systems, 35:34586–34599, 2022.
660	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timo-
661 662	thy Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. <i>arXiv</i> preprint arXiv:2306.09212, 2023a.
663	
664 665	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Halueval: A large- scale hallucination evaluation benchmark for large language models. In <i>Proceedings of the 2023</i>
666	Conjerence on Empirical Methods in Natural Language 1 locessing, pp. 0449–0404, 20250.
667	Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong
668	Wen. The dawn after the dark: An empirical study on factuality hallucination in large language models, 2024a. URL https://arxiv.org/abs/2401.03205.
669	
670	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time
670	intervention: Eliciting truthful answers from a language model. Advances in Neural Information
672	Processing Systems, 36, 20246.
674	Tianle Li, Ge Zhang, Ouv Duc Do, Xiang Yue, and Wenhu Chen, Long-context llms struggle with
675	long in-context learning. CoRR, 2024c.
676	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
677	Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following
678	models.https://github.com/tatsu-lab/alpaca_eval,52023c.
679	Chin You Lin Bouges A neckage for automatic avaluation of summarized. In Text summarization
680 681	branches out, pp. 74–81, 2004.
682	Sheng-Chieh Lin, Luvu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen-tau Yih, and
683	Xilun Chen. Flame: Factuality-aware alignment for large language models. arXiv preprint
004	arXiv:2405.01525, 2024.
600	Stephanie Lin Jacob Hilton and Owain Evans, Truthfulga: Measuring how models mimic human
000	falsehoods In Proceedings of the 60th Annual Meeting of the Association for Computational
688	Linguistics (Volume 1: Long Papers), pp. 3214–3252, 2022.
689	Nelson F Liu Kevin Lin John Hewitt Ashwin Paraniana Michala Ravilacoua Fahio Patroni and
690	Percy Liang Lost in the middle: How language models use long contexts. Transactions of the
691	Association for Computational Linguistics 12:157–173 2024a
692	19900 and 197 Comparational Englishes, 12.197 119, 2027a.
693	Tianyu Liu, Yizhe Zhang, Chris Brockett, Yi Mao, Zhifang Sui, Weizhu Chen, and William B Dolan.
694	A token-level reference-free hallucination detection benchmark for free-form text generation. In
695	Proceedings of the outh Annual Meeting of the Association for Computational Linguistics (Volume 1. Long Papers) pp. 6723–6737, 2022
696	1. Long 1 aperol, pp. 0125-0151, 2022.
697	Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke,
698	Yifan Xu, Weng Lam Tam, et al. Alignbench: Benchmarking chinese alignment of large language
699	models. arXiv preprint arXiv:2311.18743, 2023.
700	Vang Liu Jiahuan Cao, Changgu Liu, Kai Ding, and Lianuan Lin. Detects for large large
701	models: A comprehensive survey. <i>arXiv preprint arXiv:2402.18041</i> , 2024b.

702 703 704	Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer. <i>arXiv preprint arXiv:2405.16436</i> , 2024c.
705 706	I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
707 708 709	Renze Lou, Kai Zhang, and Wenpeng Yin. Large language model instruction following: A survey of progresses and challenges. <i>Computational Linguistics</i> , pp. 1–10, 2024.
710 711	Michelle M Mello and Neel Guha. Chatgpt and physicians' malpractice risk. In <i>JAMA Health Forum</i> , volume 4, pp. e231938–e231938. American Medical Association, 2023.
712 713 714 715 716	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 12076–12100, 2023.
717 718 719	Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia Tsvetkov, and Hannaneh Hajishirzi. Fine-grained hallucination detection and editing for language models. <i>arXiv preprint arXiv:2401.06855</i> , 2024.
720 721 722 723	Fedor Moiseev, Zhe Dong, Enrique Alfonseca, and Martin Jaggi. Skill: Structured knowledge infusion for large language models. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pp. 1581–1588, 2022.
724 725 726 727	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christo- pher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> , 2021.
728 729 730 731	Neel Nanda, Senthooran Rajamanoharan, János Kramár, and Rohin Shah. Fact find- ing: Attempting to reverse-engineer factual recall on the neuron level, Dec 2023. URL https://www.alignmentforum.org/posts/iGuwZTHWb6DFY3sKB/ fact-finding-attempting-to-reverse-engineer-factual-recall.
732 733 734 735	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35: 27730–27744, 2022.
736 737 738 739	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pp. 311–318, 2002.
740 741 742	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
743 744 745	N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> , 2019.
746 747 748	Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 5418–5426, 2020.
749 750 751 752	Cicero Nogueira dos Santos, Zhe Dong, Daniel Cer, John Nham, Siamak Shakeri, Jianmo Ni, and Yun-Hsuan Sung. Knowledge prompts: Injecting world knowledge into language models through soft prompts. <i>arXiv preprint arXiv:2210.04726</i> , 2022.
753 754 755	Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bow- man, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, et al. Towards understanding sycophancy in language models. In <i>The Twelfth International Conference on</i> <i>Learning Representations</i> , 2024.

756 757	Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in rlhf. <i>arXiv preprint arXiv:2310.03716</i> , 2023.
758 759	Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu Ding, Hongyang Li, Mengzhe Geng, et al. A survey of reasoning with foundation models
760 761	arXiv preprint arXiv:2312.11562, 2023.
762	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: A large-
763	scale dataset for fact extraction and verification. In 2018 Conference of the North American Chap-
764 765	ter of the Association for Computational Linguistics: Human Language Technologies, NAACL HLT 2018, pp. 809–819. Association for Computational Linguistics (ACL), 2018.
766	Katherine Tian Eric Mitchell Huaxiu Yao Christopher Manning and Chelsea Finn Fine-tuning
767 768	language models for factuality. In <i>NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following</i> , 2023.
769	
770 771	<i>of the 27th International Joint Conference on Artificial Intelligence</i> , pp. 4950–4957, 2018.
772 773	Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(11), 2008.
774	Ashish Vaswani Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez
776	Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa- tion processing systems, 30, 2017.
778	Disiis Ware Ether Charry and Darsfei Lin. Chinesefectural: A factuality handwards for chinese
779	lims 2023a
780	millo, 2020a.
781	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi
782 783	Knowledge, retrieval and domain-specificity. <i>arXiv preprint arXiv:2310.07521</i> , 2023b.
784 785	Hongmin Wang. Revisiting challenges in data-to-text generation with fact grounding. In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pp. 311–322, 2019.
786	Jason Wei, Yuazhi Wang, Dala Schuurmans, Maartan Rosma, Fai Yia, Ed Chi, Ouos V La, Danny
787 788	Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.
789	Jerry Wei, Chengrun Vang, Xinying Song, Vifeng Lu, Nathan Hu, Dustin Tran, Daiyi Peng, Ruiho
791	Liu, Da Huang, Cosmo Du, et al. Long-form factuality in large language models. <i>arXiv preprint</i>
792	arXiv:2403.18802, 2024.
793	Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenije Ye, Weilin Liu, Zhivu Mei, Guangiu Wang, Chao Yu
794	and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. In <i>Forty-first</i>
795	International Conference on Machine Learning, 2024.
796	Aiyuan Yang Bin Xiao Bingning Wang Borong Zhang Ce Bian Chao Vin Chenyu Ly Da Pan
797	Dian Wang, Dong Yan, et al. Baichuan 2: Onen large-scale language models. arXiv preprint
798	arXiv:2309.10305, 2023.
800	An Vana Dessana Vana Dimmon II.: De Zhana Demon Va Chana Zhau Chananana Li
801	All rang, Dausong rang, Dinyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengwian Li Daviheng Liu Fei Huang et al. Owen? technical report arViv preprint
802	arXiv:2407.10671. 2024.
803	
804	Y1 Yang, Wen-tau Y1h, and Christopher Meek. W1k1qa: A challenge dataset for open-domain ques- tion answering. In <i>Proceedings of the 2015 conference on curricities with a data in actual language</i>
805	non answering. In Proceedings of the 2015 conference on empirical methods in natural language processing on 2013–2018 2015
806	processing, pp. 2015–2010, 2015.
807	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,
808	and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question
809	Processing, pp. 2369–2380, 2018.

810 811 812	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large language models know what they don't know? In <i>The 61st Annual Meeting Of The Association For Computational Linguistics</i> , 2023.
814 815	Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. Alleviating hallucinations of large language models through induced hallucinations. <i>arXiv preprint arXiv:2312.15710</i> , 2023a.
816 817 818	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. Siren's song in the ai ocean: a survey on hallucination in large language models. <i>arXiv preprint arXiv:2309.01219</i> , 2023b.
819 820 821 822	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. <i>arXiv</i> preprint arXiv:2303.18223, 2023.
823 824 825	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
826 827 828	
829 830 831	
832 833	
834 835 836	
837 838	
839 840 841	
842 843 844	
845 846	
847 848 849	
850 851	
852 853 854	
855 856 857	
858 859	
860 861 862	
863	

864 APPENDIX

866

868

879

882 883

885 886

A LOW-QUALITY DATA IN PRE-GENERATION

In pre-generation experiment, we randomly sampled 3,000 encyclopedia entries and directly used
simple prompts to generate approximately 9,000 questions. After manually evaluating the questions
and answers, we found four types of low-quality data: 1) The related knowledge is too long-tailed;
2) The answer to the question is not unique; 3) The answer generated by GPT4 is incorrect; 4) The
question is not self-contained.

In this section, we will provide a detailed introduction to these typical low-quality data. In the
subsequent generation process, we focused on how to avoid the generation of such low-quality data
and formulated a better pipeline. The basic introductions and solutions of four typical types of
low-quality data are shown in table 7.

Table 7: Types of low-quality data found in pre-generation and solutions.

Types of low-quality data	Causes or features	Solutions
The related knowledge is too long-tailed	The encyclopedia entry is relatively unknown and useless	Filter long-tailed entries based on encyclopedia view counts
The correct answer to the question is not unique	The knowledge in encyclopedia entry's description can- not cover all the facts in the world	Use carefully designed prompts for assistance when gen- erating questions, and further filter after generating ques- tions
The answer generated by GPT4 is incorrect	The encyclopedia entry's description is too long, which leads to a difficulty in the understanding and instruct fol- lowing	Filter based on the length of the encyclopedia entry's de- scription
	The answer to the question depends on the time limitation (a time-sensitive question)	Use carefully designed prompts for assistance when gen- erating questions, and further filter after generating ques- tions
	The encyclopedia entry's description is difficult for model to understand	Low frequency, no additional processing at present
The question is not self-contained	The question is not accurate and specific enough	Use carefully designed prompts for assistance when gen- erating questions, and further filter after generating ques- tions
	The encyclopedia entry's description itself is not self- contained	Low frequency in entries with high view counts. Filter based on encyclopedia views counts

896 897

898

899 900

901

The related knowledge is too long-tailed Some entries' subjects are relatively unknown, and the related knowledge appears less on the internet since it is not important for most of people. Questions on those subjects have little value, and are too hard to be answered correctly.

The correct answer to the question is not unique The answers to some questions are not unique and may extend beyond the knowledge provided in encyclopedia entries' descriptions. Additionally, some terms in questions involve subjective judgments hence questions will have more possible and reasonable answers.

 The answer generated by GPT4 is incorrect GPT4 is applied to output both question and standard answer based on the description of encyclopedia entry. As the description's length increases, it will challenge model's ability to follow instruction and catch the precise knowledge of the description. Sometimes the description itself is difficult to understand, such as poetry and classical Chinese article, and GPT4 could provide incorrect answers or even semantically incoherent questions facing them. Another situation is that some encyclopedia entries do not emphasize the time limitation of information, and for some rapidly changing data, the knowledge could be outdated.

The question is not self-contained During benchmarking, only the question is provided for the
 model without encyclopedia content. Therefore, the question should be understandable, complete,
 and cannot contain vague pronouns or nouns with multiple interpretations such as abbreviations
 and names without clear context, unless in the vast majority of cases, the word has the same ref erence or meaning. Sometimes low-quality encyclopedia entries' descriptions themselves are not self-contained.

918 B MORE DETAILS ABOUT FACTUALBENCH 919

B.1 PROMPTS

920

921 922

923

924

926

927

In this section, we will provide all prompts related to FactualBench, including generation part and evaluation part.

- 925 B.1.1 GENERATION
 - **Question generation** (English translation is for reference only):
- 928
 929 我将提供给你一个对象和相关的参考文档,请针对对象提出最多{提问个数:3}个事实性
 930 问题。要求每个问题都具有唯一且准确的答案,避免答案模糊或存在争议,避免涉及主观
 931 判断的问题和时效性问题,要求答案可以在参考文档中直接找到。要求提问的问题表达清
 932 断。问题中的名词指长出现,不需要优格会考文档中直接找到。要求提问的问题表达清
- 晰,问题中的名词指代明确,不需要依赖参考文档即可理解问题内容。对每个问题,给出1 933 个标准答案和3个具有干扰性的错误答案。 934 下面是两个例子: 935 936 【对象】: {示例对象1} 937 【参考文档】: {关于示例对象1的百科内容简介} 938 【问题1】: {针对示例对象1提出的示例问题1} 939 940 【标准答案】: {示例问题1标准答案} 941 【错误答案1】: {示例问题1错误答案1} 942 943 【错误答案2】: {示例问题1错误答案2} 944 945 【错误答案3】: {示例问题1错误答案3} 946 【对象】: {示例对象2} 947 948 【参考文档】: {关于示例对象2的百科内容简介} 949 【问题1】: {针对示例对象2提出的示例问题2} 950 951 【标准答案】: {示例问题2标准答案} 952 【错误答案1】: {示例问题2错误答案1} 953 954 【错误答案2】: {示例问题2错误答案2} 955 【错误答案3】: {示例问题2错误答案3} 956 对于以下的对象和参考文档、使用同样的格式生成问题、答案。 957 958 【对象】: {对象: 百科词条对象} 959
 - 【参考文档】: {文档: 百科简介}

968

- 969 Here are two examples:
- 970 [Object]: {Example Object 1} 971

[Reference Description]: {Brief introduction to Example Object 1}

⁹⁶⁰ 961 962

^{I will provide you with an object and related reference description. Please generate up to {Question number: 3} factual questions about the object. Each question should have a unique and accurate answer, avoiding vague or contentious answers, subjective judgments, and time-sensitive. The answers should be directly found in the reference description. The questions should be clearly expressed, with unambiguous noun references, and should not rely on the reference description for understanding. For each question, provide one correct answer and three misleading incorrect answers.}

- 972 [Question 1]: {Example question 1 related to Example Object 1}
- 974 [Correct Answer]: {Standard answer to Example question 1}
- 975 [Incorrect Answer 1]: {Incorrect answer 1 to Example question 1} 976
- [Incorrect Answer 2]: {Incorrect answer 2 to Example question 1}
- 978 [Incorrect Answer 3]: {Incorrect answer 3 to Example question 1}
- 979 [Object]: {Example Object 2}
- 981 [Reference Document]: {Brief introduction to Example Object 2}
- 982 [Question 1]: {Example question 2 related to Example Object 2}
- [Correct Answer]: {Standard answer to Example question 2}
- 985 [Incorrect Answer 1]: {Incorrect answer 1 to Example question 2}
- 986 987 [Incorrect Answer 2]: {Incorrect answer 2 to Example question 2}
- 988 [Incorrect Answer 3]: {Incorrect answer 3 to Example question 2}
- For the following object and reference description, generate questions and answers in the same format.
- [Object]: {Object: Encyclopedia Entry Object}
- 993 [Reference Description]: {Description: Encyclopedia Description}

994 995 996

997 998 **Question filtering** (English translation is for reference only):

999 User:你是一个评估专家,下面需要你对一个问题的质量进行判断。

- 1000 我会给你一个事实性知识问答问题,你需要从以下几个方面分析这个问题,最终回答问题是1001 【优质】还是【非优质】。
- 1002 如果这个问题内存在代词指代不清,或无法明确理解问题含义,请回复【非优质】。
- 1004 如果问题的答案不唯一,请回复【非优质】。
- 1005 如果问题是时效性问题,且没有给出具体的背景时间点,请回复【非优质】。 1006
- 1007 如果问题没有以上情况,请回复【优质】。
- 1008 请一步步思考,并在最后给出你的判断:【优质】或【非优质】。注意将你的最终判断写 1009 在中括号【】中!
- **1010** Assistant:明白了,我会按照你的要求和规则进行判断
- 1012 User:问题是:
- 1013 {待评价问题}
- 1014 1015 请给出你的判断:
- 1016

User: You are an evaluation expert, and you need to assess the quality of a question.

I will provide you with a factual knowledge question, and you need to analyze the question from thefollowing aspects to determine whether the question is of [High Quality] or [Low Quality].

If the question contains unclear pronoun references or cannot be clearly understood, please respond with [Low Quality].

- If the answer to the question is not unique, please respond with [Low Quality].
- 1025

1026 If the question is time-sensitive and does not provide a specific time limitation, please respond with 1027 [Low Quality]. 1028 1029 If none of the above situations apply, please respond with [High Quality]. 1030 Please think through the question step by step and give your final judgment as [High Quality] or 1031 [Low Quality]. Be sure to put your final judgment in square brackets []! 1032 Assistant: Understood. I will follow your requirements and rules for evaluation. 1033 1034 User: The question is: 1035 {Question to be evaluated} 1036 1037 Please provide your judgment: 1038 1039 1040 1041 **B.1.2** EVALUATION 1042 1043 1044 Answer evaluation (English translation is for reference only): 1045 1046 给定一个问题以及对应的参考答案,根据参考答案和你掌握的知识,对候选答案是否回答 1047 正确进行评价。请注意,问题可能不具有唯一答案,此时只要候选答案言之成理即可。如 1048 果候选答案符合参考答案或言之成理,请回答【正确】;如果候选答案与参考答案矛盾或 1049 没有回答问题,请回答【错误】,并给出你的分析过程。下面是五个例子: 1050 【问题】: 百川智能创始人王小川在什么时间与茹立云联合创立了该公司? 1051 1052 【参考答案】: 百川智能创始人王小川于2023年4月10日与茹立云联合创立了该公司。 【候选答案】: 王小川与茹立云于2023年4月共同创立了百川智能公司。 1054 1055 【评价】: 根据参考答案, 百川智能于2023年4月10日创立, 候选答案认为是2023年4月创 1056 立,符合参考答案。【正确】 1057 【问题】:《采桑子·清明后三日作》是哪位诗人创作的? 1058 【参考答案】:《采桑子·清明后三日作》是诗人龙榆生创作的。 1059 【候选答案】:《采桑子·清明后三日作》是清代词人蒋春霖创作的一首词。 1061 【评价】:根据参考答案,《采桑子·清明后三日作》是由诗人龙榆生创作,候选答案认为 1062 是蒋春霖创作,与参考答案矛盾。【错误】 1063 1064 【问题】: 李白的代表作有哪些? 【参考答案】: 李白的代表作有《望庐山瀑布》《行路难》 《蜀道难》 《将进酒》 《早发 白帝城》《黄鹤楼送孟浩然之广陵》等。 1067 【候选答案】: 李白的代表作有《将进酒》《静夜思》《庐山谣》 《早发白帝城》 《赠汪 1068 伦》 《望庐山瀑布》《行路难》《夜泊牛渚怀古》 《登金陵凤凰台》 《送友人》等。 1069 【评价】: 李白有许多代表作, 答案不唯一, 候选答案中的诗的确均为李白所写, 言之成 1071 理。【正确】 【问题】: 哈蒂·温斯顿的主要作品有哪些? 1073 【参考答案】:哈蒂·温斯顿的主要作品有《灵书妙探第一季》。 1074 1075 【候选答案】:哈蒂·温斯顿(Hedy Lamarr)的主要作品有《Ecstasy》 (1933年), 1076 《Algiers》(1938年), 《Samson and Delilah》 (1949年)等。 1077 【评价】: 哈蒂·温斯顿有许多作品, 答案不唯一, 但候选答案中的作品不是哈蒂·温斯顿的 1078 作品。【错误】 1079 【问题】: 吴之番在哪次战斗中牺牲的?

1094

1080 【参考答案】: 吴之番在清顺治二年八月二十六日的战斗中牺牲, 这是嘉定三屠的一部 1081 分。 1082 【候选答案】:对不起,我找不到关于"吴之番"的相关牺牲信息。这可能是因为您提供的 1083 信息有误或者该人物并不存在。 1084 【评价】:根据参考答案,吴之番在顺治二年八月二十六日的战斗中牺牲,候选答案没有 回答问题。【错误】 1086 1087 下面是你需要评价的内容,请使用同样的格式给出评价。 1088 【问题】: {问题} 1089 1090 【参考答案】:{参考答案} 1091 【候选答案】:{候选答案} 1092 1093

Given a question and its corresponding standard answer, evaluate whether the candidate answer correctly addresses the question based on the standard answer and your knowledge. Please note that the question may not have only one unique answer; in such cases, as long as the candidat answer is reasonable, it is acceptable. If the candidate answer aligns with the reference answer or is reasonable, please respond with [Correct]; if the candidate answer contradicts the reference answer or refuses to answer the question, please respond with [Incorrect] and provide your analysis. Here are five examples:

- [Question]: When did Wang Xiaochuan, the founder of Baichuan Inc., co-found the company with Ru Liyun?
- [1104 [Standard Answer]: Wang Xiaochuan co-founded Baichuan Inc. with Ru Liyun on April 10, 2023.
- [Candidate Answer]: Wang Xiaochuan and Ru Liyun co-founded Baichuan Inc. in April 2023.
- [Evaluation]: According to the standard answer, Baichuan Inc. was founded on April 10, 2023.
 The candidate answer states it was founded in April 2023, which aligns with the reference answer.
 [Correct]
- 1110 [Question]: Which poet created "Cai Sang Zi · Qing Ming Hou San Ri Zuo"?
- [Standard Answer]: "Cai Sang Zi · Qing Ming Hou San Ri Zuo" was created by the poet Long Yusheng.
- [Candidate Answer]: "Cai Sang Zi · Qing Ming Hou San Ri Zuo" was created by the Qing Dynasty poet Jiang Chunlin.
- [Evaluation]: According to the reference answer, "Cai Sang Zi · Qing Ming Hou San Ri Zuo" was
 created by Long Yusheng, while the candidate answer claims it was created by Jiang Chunlin, which
 contradicts the reference answer. [Incorrect]
- ¹¹¹⁹ [Question]: What are the representative works of Li Bai?
- [Standard Answer]: Li Bai's representative works include "Wang Lu Shan Pu Bu", "Xing Lu Nan",
 "Shu Dao Nan", "Qiang Jin Jiu", "Zao Fa Bai Di Cheng", and "Huang He Lou Song Meng Hao Ran Zhi Guang Ling", etc.
- [Candidate Answer]: Li Bai's representative works include "Qiang Jin Jiu", "Jing Ye Si", "Lu Shan Yao", "Zao Fa Bai Di Cheng", "Zeng Wang Lun", "Wang Lu Shan Pu Bu", "Xing Lu Nan", "Ye Bo Niu Zhu Huai Gu", "Deng Jin Ling Feng Huang Tai", and "Song You Ren", etc.
- [Evaluation]: Li Bai has many representative works, and the answer is not unique. The poems listed
 in the candidate answer are indeed all written by Li Bai, which is reasonable. [Correct]
- 1129 [Question]: What are the main works of Hattie Winston?
- 1131 [Standard Answer]: Hattie Winston's main work is "Castle" (Season one).
- [Candidate Answer]: Hedy Lamarr's main works include "Ecstasy" (1933), "Algiers" (1938), and
 "Samson and Delilah" (1949), etc.

1134 [Evaluation]: Hattie Winston has many works, and the answer is not unique. However, the works 1135 listed in the candidate answer are not by Hattie Winston. [Incorrect] 1136 [Question]: In which battle did Wu Zhifan sacrifice? 1137 1138 [Standard Answer]: Wu Zhifan was sacrificed in the battle on August 26, the second year of the 1139 Shunzhi reign, which was part of the Jiadin Santu. 1140 [Candidate Answer]: Sorry, I cannot find any information related to Wu Zhifan's sacrifice. This may 1141 be due to incorrect information you provided or because this person does not exist. 1142 [Evaluation]: According to the standard answer, Wu Zhifan was sacrificed in the battle on August 1143 26, the second year of the Shunzhi reign, but the candidate answer did not answer the question. 1144 [Incorrect] 1145 1146 Here is the content you need to evaluate, and please use the same format to provide your evaluation. 1147 [Question]: {Question} 1148 [Standard Answer]: {Standard Answer} 1149 1150 [Candidate Answer]: {Candidate Answer} 1151 1152 1153 **B.2** MORE EXAMPLE 1154 1155 We show one example of each domain from FactualBench in Table 8. 1156 1157 1158 C DETAILED SETTINGS 1159 1160 In this section, we will introduce the settings of our experiments, including the evlaution and training 1161 parts. 1162 1163 C.1 EVALUATION 1164 1165 We benchmark 14 LLMs on our FactualBench: Baichuan1 (closed), Baichuan2 (open), Qwen1.5-7B-Instruct (open), Qwen2-7B-Instruct (open), Llama-3-8B-Instruct (open), Baichuan3 (closed), 1166 Yi-34B-Chat (open), Command-R-35B (open), Llama-3-70B-Instruct (open), Qwen2-72B-Instruct 1167 (open), Baichuan4 (closed), Command-R-plus-104B (open), DeepSeek-v2-0628 MoE-236B (open), 1168 GPT4-0125-preview (closed), among which DeepSeek and GPT4 are queried from api and others 1169 are inferenced locally. 1170 We prioritize the chat/instruct version of the model and use the providing recommended generation 1171 config and code on huggingface² to generate responses. We set max_new_tokens or max_length 1172 configuration large enough to ensure that models can complete their normal responses. 1173 1174 C.2 EXPERIMENTS 1175 1176 We choose 6 other open source benchmarks to evaluate model's enhancement comprehensively. We 1177 generate answer zero-shot inputing the questions or instructions in default generation config. For 1178 model-base evaluation process, we all choose GPT4-0125-preview as evaluator. 1179 TruthfulQA Lin et al. (2022) is an English benchmark to measure whether a language model is 1180 truthful in generating answers. It contains 817 questions covering 38 domains. The questions are 1181 designed to cause imitative falsehoods, a false may due to wrong training objective like fake knowl-1182 edge in training data. We use the generative part of TruthfulQA and use GPT4 to evaluate the answer 1183 (provide reference correct and incorrect answers when judging). 1184 1185 HalluQA (Cheng et al., 2023) is a benchmark to measure hallucination phenomenon in Chinese 1186 LLM. It contains 450 meticulously designed adversarial questions covering diverse domains to test 1187

²https://huggingface.co

Question	Standard answer	Wrong answer 1	Wrong answer 2	Wrong answer 3	Domain
韩国电影《人狼》是由哪位导演执导 的?	电影《人狼》是由金知云执导的。	电影《人狼》是由姜栋元执导的。	电影《人狼》是由韩孝周执导的。	电影《人狼》是由郑雨盛执导的。	影视娱乐
河北师范大学最早起源于哪两所学校?	河北师范大学最早起源于顺天府学堂和 北洋女师范学堂。	河北师范大学最早起源于河北师范学院 和河北教育学院。	河北师范大学最早起源于河北职业技术 师范学院和汇华学院。	河北师范大学最早起源于北京大学和清 华大学。	教育培养
苯丙氨酸的化学式是什么?	苯丙氨酸的化学式是C9H11NO2。	苯丙氨酸的化学式是C8H11NO2。	苯丙氨酸的化学式是C9H10NO2。	苯丙氨酸的化学式是C9H11NO3。	数理化生
谥号是在什么时期开始的?	谥号始于西周-	谥号始于东周-	谥号始于秦朝-	谥号始于汉朝-	历史国等
中国电影"第六代导演"之一王小帅的电 影处女作是什么?	王小帅的电影处女作是《冬春的日 子》,	王小帅的电影处女作是《扁担姑娘》.	王小帅的电影处女作是《十七岁的单 车》。	王小帅的电影处女作是《青红》 -	人物百利
法律关系的构成要素有哪些?	法律关系的构成要素有三项:法律关系 主体,法律关系内容,法律关系客体。	法律关系的构成要素有三项:法律关系 主体,法律关系形式,法律关系客体。	法律关系的构成要素有三项:法律关系 主体,法律关系内容,法律关系方式。	法律关系的构成要素有三项:法律关系 主体,法律关系内容,法律关系目标。	政治法律
国家金融监督管理总局是在哪一年揭牌 的?	国家金融监督管理总局是在2023年揭牌 的-	国家金融监督管理总局是在2022年揭牌 的。	国家金融监督管理总局是在2021年揭牌 的。	国家金融监督管理总局是在2020年揭牌 的。	经济管理
MemCache是由谁开发的?	MemCache是由LiveJournal的Brad Fitz- patrick开发的。	MemCache是由Facebook的Mark Zucker- berg开发的。	MemCache是由Google的Larry Page开发的.	MemCache是由Microsoft的Bill Gates开发的。	计算机科
瑞舒伐他汀的主要作用部位是哪里?	瑞舒伐他汀的主要作用部位是肝。	瑞舒伐他汀的主要作用部位是心脏。	瑞舒伐他汀的主要作用部位是肾脏。	瑞舒伐他汀的主要作用部位是胃。	医学
"枫丹白露"这个名字的原义是什么?	"枫丹白露"的法文原义为"美丽的泉 水"-	"枫丹白露"的法文原义为"宏伟的宫 殿"-	"枫丹白露"的法文原义为"狩猎的行 宫"。	"枫丹白露"的法文原义为"古老的城堡"。	社会人文
竹笋原产于哪里?	竹笋原产于中国。	竹笋原产于日本。	竹笋原产于印度。	竹笋原产于泰国.	农林牧道
更新世是由哪位地质学家创用的?	更新世是由英国地质学家菜伊尔创用 的。	更新世是由英国地质学家福布斯创用 的。	更新世是由美国地质学家莱伊尔创用 的,	更新世是由中国地质学家莱伊尔创用 的-	天文地理
新奥尔良鹈鹕队在哪一年正式宣布球队 改名为鹈鹕队?	新奥尔良鹈鹕队在2013年正式宣布球队 改名为鹈鹕队。	新奥尔良鹈鹕队在2012年正式宣布球队 改名为鹈鹕队。	新奥尔良鹈鹕队在2014年正式宣布球队 改名为鹈鹕队。	新奥尔良鹈鹕队在2015年正式宣布球队 改名为鹈鹕队。	运动旅游
宾利汽车公司是在哪一年创办的?	宾利汽车公司是在1919年创办的。	宾利汽车公司是在1920年创办的。	宾利汽车公司是在1918年创办的。	宾利汽车公司是在1921年创办的。	数码汽车
隔离开关主要用于什么?	隔离开关主要用于隔离电源、倒闸操 作、用以连通和切断小电流电路。	隔离开关主要用于调节电压。	隔离开关主要用于转换电流。	隔离开关主要用于存储电能。	工业工程
鸦片战争是在哪一年开始的?	鸦片战争是在1840年开始的-	鸦片战争是在1842年开始的。	鸦片战争是在1839年开始的。	鸦片战争是在1841年开始的。	军武战争
买了佛冷这个词是来源于哪首歌曲?	买了佛冷这个词是来源于歌曲《I Love Poland》	买了佛冷这个词是来源于歌曲《I Love China》。	买了佛冷这个词是来源于歌曲《I Love America》 -	买了佛冷这个词是来源于歌曲《I Love England》	网词网梗
苏荷酒吧是在哪一年诞生的?	苏荷酒吧是在2003年诞生的。	苏荷酒吧是在2000年诞生的。	苏荷酒吧是在2005年诞生的。	苏荷酒吧是在2010年诞生的-	工作生活
视觉识别系统VI是什么的缩写?	視觉识别系统是Visual Identity的缩写。	视觉识别系统是Visual Information的缩写。	视觉识别系统是Visual Interface的缩写。	视觉识别系统是Visual Interaction的缩写。	高新科技
风水业内公认的"龙脉之源"是哪里?	风水业内公认的"龙脉之源"是昆仑山-	风水业内公认的"龙脉之源"是长江。	风水业内公认的"龙脉之源"是黄河 -	风水业内公认的"龙脉之源"是太湖。	信仰文化
Who directed the Korean movie 'Inrang'?	The movie 'Inrang' is directed by Kim Jee- woon.	The movie 'Inrang' is directed by Kang Dong Won.	The movie 'Inrang' is directed by Han Hy- oJoo.	The movie 'Inrang' is directed by Jung Woo Sung.	film&entertair
Which two schools did Hebei Normal Uni- versity first originate from?	Hebei Normal University originated from Shuntianfu Official School and Beiyang Women's Normal School.	Hebei Normal University originated from Hebei Normal Institute and Hebei Institute of Education.	Hebei Normal University originated from Hebei Vocational and Technical Normal College and Huihua College.	Hebei Normal University originated from Peking University and Tsinghua Univer- sity.	education&tra
What is the chemical formula for pheny- lalanine?	The chemical formula for phenylalanine is C9H11NO2.	The chemical formula for phenylalanine is C8H11NO2.	The chemical formula for phenylalanine is C9H10NO2.	The chemical formula for phenylalanine is C9H11NO3.	physics, chemistry, mathe
When did posthumous titles begin?	The posthumous title began in the Western Zhou Dynasty.	The posthumous title began in the Eastern Zhou Dynasty.	The posthumous title began in the Qin Dy- nasty.	The posthumous title began in the Han Dy- nasty.	history&tradition
What is the debut film of Wang Xiaoshuai, one of the "sixth generation directors" of Chinace giagene?	Wang Xiaoshuai's debut film is 'THE DAYS'.	Wang Xiaoshuai's debut film is 'So Close to Paradise'.	Wang Xiaoshuai's debut film is 'Beijing Bicycle'.	Wang Xiaoshuai's debut film is 'Shanghai Dreams'.	biography
What are the constituent elements of legal	There are three elements that make up a le-	There are three elements that make up a	There are three elements that make up a le-	There are three elements that make up a le-	politics&la
relationships?	gal relationship: the subject of the legal re-	legal relationship: the subject of the legal	gal relationship: the subject of the legal re-	gal relationship: the subject of the legal re-	
	lationship, the content of the legal relation- ship, and the object of the legal relation- ship.	relationship, the form of the legal relation- ship, and the object of the legal relation- ship.	lationship, the content of the legal relation- ship, and the method of the legal relation- ship.	lationship, the content of the legal relation- ship, and the objective of the legal relation- ship.	
In which year was the Chinese National	The Chinese National Financial Supervi-	The Chinese National Financial Supervi-	The Chinese National Financial Supervi-	The Chinese National Financial Supervi-	economics&man:
Financial Supervisory Administration un- veiled?	sory Administration was unveiled in 2023.	sory Administration was unveiled in 2022.	sory Administration was unveiled in 2021.	sory Administration was unveiled in 2020.	
Who developed MemCache?	MemCache was developed by Brad Fitz- patrick from LiveJournal.	MemCache was developed by Mark Zuckerberg from Facebook.	MemCache was developed by Larry Page from Google.	MemCache was developed by Bill Gates from Microsoft.	computer sci
What is the main site of action of rosuvas- tatin?	The main site of action of rosuvastatin is the liver.	The main site of action of rosuvastatin is the heart.	The main site of action of rosuvastatin is the kidney.	The main site of action of rosuvastatin is the stomache.	medical
What is the original meaning of 'Fontainebleau'?	The original French meaning of "Fontainebleau" is "beautiful spring water".	The original French meaning of "Fontainebleau" is "magnificent palace".	The original French meaning of "Fontainebleau" is "hunting palace".	The original French meaning of "Fontainebleau" is "ancient castle".	sociology&hun
Where do bamboo shoots originate from?	Bamboo shoots originate from China.	Bamboo shoots originate from Japan.	Bamboo shoots originate from India.	Bamboo shoots originate from Thailand.	agriculture, forestry, fisherie
Which geologist named the Pleistocene epoch?	The Pleistocene was named by British ge- ologist Lyell.	The Pleistocene was named by British ge- ologist Forbes.	The Pleistocene was named by American geologist Lyell.	The Pleistocene was named by Chinese ge- ologist Lyell.	astronomy&geo
In which year did the New Orleans Peli- cans officially announce their name change	The New Orleans Pelicans officially an- nounced their name change to the Pelicans	The New Orleans Pelicans officially an- nounced their name change to the Pelicans	The New Orleans Pelicans officially an- nounced their name change to the Pelicans	The New Orleans Pelicans officially an- nounced their name change to the Pelicans	sports&tour
to the Pelicans? In which year was BentleyMotors Limited	in 2013. BentleyMotors Limited was founded in	in 2012. BentleyMotors Limited was founded in	in 2014. BentleyMotors Limited was founded in	in 2015. BentleyMotors Limited was founded in	digital&autor
founded? What is the main use of disconnectors?	1919. Disconnectors are mainly used for isolat-	1920. Disconnectors are mainly used to regulate	1918. Disconnectors are mainly used to convert	1921. Disconnectors are mainly used for storing	industrial engir
	ing power sources, switching operations, and connecting and disconnecting small current circuits.	voltage.	current.	electrical energy.	
In which year did the Opium War begin?	The Opium War begin in 1840.	The Opium War begin in 1842.	The Opium War begin in 1839.	The Opium War begin in 1841.	military&w
What song does the meme 'Mai Le Fo Leng' come from?	The meme 'Mai Le Fo Leng' comes from "I love Poland"	The meme 'Mai Le Fo Leng' comes from "I love China"	The meme 'Mai Le Fo Leng' comes from "I love America"	The meme 'Mai Le Fo Leng' comes from "I love England"	slang&mer
In which year was Soho Bar founded?	Soho Bar was founded in 2003.	Soho Bar was founded in 2000.	Soho Bar was founded in 2005.	Soho Bar was founded in 2010.	work&lif
What word is VI(a Vision System) abbre-	VI abbreviation for Visual Identity.	VI abbreviation for Visual Information.	VI abbreviation for Visual Interface.	VI abbreviation for Visual Interface.	high technol
viation for?					

Table 8: More examples from FactualBench. English translation is for reference only.

1224

1188

1225

models imitative falsehoods and knowledge. Still, we use the generative part and its official prompt to evaluate the answer.

1228 CMMLU (Li et al., 2023a) is a Chinese multiple-choice benchmark similar to MMLU Hendrycks
 et al. (2020), comprising 67 topics and massive questions. We use the official script and code to
 evaluate model's accuracy by logits output.

HaluEval (Li et al., 2023b) is a large collection of generated and human-annotated English hallucinated samples for evaluating the performance of LLMs in recognizing hallucination. It's a discriminative task require test model to judge whether an answer contains hallucination or not. We use the official prompt form to query, and only use 10000 samples from QA part. The evaluation is based on string matching (e.g. "Yes" or "No"), and we have added more matching patterns. If the answer doesn't match any pattern, it will be judge as a wrong answer.

Alignbench (Liu et al., 2023) is a Chinese benchmark for evaluating LLMs' alignment. It contains 683 instructions on 8 different fundamental abilities, such as writting, reasoning, role-play, etc. We use its official prompt format to evaluate model's answer in a model-base way.

AlpacaEval (Li et al., 2023c) is a benchmark based on the AlpacaFarm (Dubois et al., 2024) evaluation set, which tests model's instruction following ability. It contains 805 samples on different instructions, and uses winning rate of the answer against a base model as metric. It has been used to indicate model's helpfulness in previous works (Lin et al., 2024). In our work, the model before training is selected as the based model.

Together with our benchmark FactualBench, we can evaluate model's factuality from generative, multiple-choice, and detective dimensions, as well as Chinese and English language. Besides, we use Alignbench to measure a model's fundamental abilities and AlpacaEval for helpfulness.

1249 1250 C.3 TRAINING

We conduct training on 32 H800-80G NVIDIA GPU using AdamW optimizer (Loshchilov, 2017).
Learning rate is set to be 2e-6 for SFT and 1e-6 for DPO. We use linear scheduler or cosine scheduler with warming up. When fusing DPO with SFT, the weight ratio of DPO loss and SFT loss is 10:1.
And we only train one epoch on the tuning set for each method and setting.

1255 1256

1257

D DETAILED BENCHMARK RESULTS

We present the performances of 14 LLMs on our FactualBench using a heatmap in Figure 8. The first column represents the overall accuracy of the model and the last line shows the average accuracy of all 14 models. We arrange domains from left to right in descending order of the average accuracy. Each domain is represent by its first 5 letters.





From Table 9, It can be observed that as models' parameters increase, the accuracy shows an upward trend, while models proficient in Chinese have a significant better performance compared to models proficient in English, which is in line with expectation. Besides, we have two more observations:
1) The same model has significantly different performances on different domains; 2) Different models share a consistency in relative ability of different domains, that is, most models' top (bottom) 5 accuracy domains are the same, and no domain exists in both top 5 and bottom 5 domains set at the same time.

1297					
1298	Model	Proficient lang.	Acc.	Top5 acc. domains(high \rightarrow low)	Bottom5 acc. domains(low \rightarrow high)
1299	Baichuan1	CN	48.24	medic, compu, high , indus, socio	film&, biogr, educa, sport, astro
1300	Baichuan2	CN	55.37	medic, physi, digit, high , compu	film&, biogr, educa, sport, milit
1301	Qwen1.5-7B Owen2-7B	CN CN	48.87 56 27	medic, compu, agric, physi, high medic, high, physi, indus, compu	film&, biogr, sport, educa, milit film& biogr educa sport milit
1302	Llama-3-8B	EN	39.11	high, compu, medic, physi, digit	film&, biogr, educa, histo, sport
1303	Baichuan3	CN	67.50	medic, physi, compu, socio, indus	sport, biogr, educa, film&, polit
1304	Yi-34B	CN	67.30	medic, compu, high , physi, socio	film&, biogr, educa, sport, histo
1004	Command-R 35B	EN	54.30	medic, compu, high , physi, indus	film&, biogr, educa, sport, histo
1305	Llama-3-70B	EN	49.65	medic, compu, high , digit, physi	film&, biogr, slang, educa, sport
1306	Qwen2-72B	CN	73.71	medic, compu, physi, high , agric	film&, biogr, educ, sport, histo
1307	Baichuan4	CN	75.07	medic, compu, physi, digit, high	film&, educa, sport, histo, biogr
1007	Command-R+ 104B	EN	60.17	medic, physi, compu, high , socio	film&, biogr, educa, sport, polic
1308	DeepSeek-v2-0628 MoE-236B	CN	75.62	medic, physi, physi, digit, high	educa, sport, biogr, film&, polit
1309	GPT4-0125-preview	EN	65.71	medic, physi, compu, relig, high	film&, biogr, educa, histo, sport

Table 9: 14 LLMs performances on FactualBench.

1310 1311

1296

1312 We demonstrate the phenomenon's occurance is due to two factors. The type of knowledge needed 1313 in different domains, and the proportion of data from different domains in the training data, which 1314 make domains' task difficulty and LLMs' mastery of domains' knowledge varies. We selected four 1315 domains with the poorest performance and four domains with the best performance almost in all 14 models, utilizing all-MiniLM-L6-v2³ from Sentence Transformer (Reimers, 2019) to extract the 1316 features and use t-SNE (Van der Maaten & Hinton, 2008) to reduce the features down to two, visu-1317 alizing in Figure 9. Different domains questions have significantly different features, comparing the 1318 best domains (mainly center at below) and the poorest domains (mainly center at above), indicating 1319 the questions distribution varies. 1320



³https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

1350 E MUTUAL NEAREST-NEIGHBOR METRIC

For two models with representations f, g, the mutual k-nearest neighbor metric measures the average overlap of their respective nearest neighbor sets (Huh et al., 2024). According to the definition (Huh et al., 2024), for each sample $x_i \sim \mathcal{X}, i = 1, 2, ..., b$, two models extract features $\phi_i = f(x_i)$ and $\psi_i = g(x_i)$. The collection of these features are denoted as $\Phi = \{\phi_1, \phi_2, ..., \phi_b\}$ and $\Psi = \{\psi_1, \psi_2, ..., \psi_b\}$. Then we compute the respective nearest neighbor sets $S(\phi_i)$ and $S(\psi_i)$ for each x_i under representations f and g.

$$d_{knn}(\phi_i, \Phi \backslash \phi_i) = S(\phi_i) \tag{4}$$

1358

1361

1362

1367

$$d_{knn}(\psi_i, \Psi \backslash \psi_i) = S(\psi_i) \tag{5}$$

where d_{knn} returns the set of indices of its k-nearest neighbors. Then we measure its average intersection via

$$m_{\rm NN}(\phi_i, \psi_i) = \frac{1}{k} |S(\phi_i) \cap S(\psi_i)| \tag{6}$$

1363 where $|\cdot|$ denotes the size of the intersection. We use Euclidean distance to calculate distance 1364 between features, and set b = 200 (if the size of dataset is less than 200, we sample all data 1365 from the dataset), k = 10 in our work. The alignment of two representations is measured by 1366 $\frac{1}{b} \sum_{i=1}^{b} m_{NN}(\phi_i, \psi_i)$.

1368 F DETAILED EXPERIMENT RESULTS

1370 F.1 ON FACTUALBENCH 1371

Baichuan1 and its related training versions' performances on FactualBench are shown in Figure 10.
 Qwen2-7B-Instruct and its related training versions' performances on FactualBench are shown in Figure 11.



Figure 10: Baichuan1 and its related training versions' performances on FactualBench.

1398 F.2 ON ALIGNBENCH

Baichuan1 and its related training versions' performances on AlignBench are shown in Table 10.
Qwen2-7B-Instruct and its related training versions' performances on AlignBench are shown in Table 11.

