# KG-FPQ: Evaluating Factuality Hallucination in LLMs with Knowledge Graph-based False Premise Questions

**Yanxu Zhu[1]**, **Jinlin Xiao[1]**, **Yuhang Wang[1]**, **Jitao Sang[1,2]** [*]
[1]Beijing Key Lab of Traffic Data Analysis and Mining, Beijing Jiaotong University
[2]Peng Cheng Lab,
**Correspondence:** yanxuzhu@bjtu.edu.cn, jinlinxiao@bjtu.edu.cn, yhangwang@bjtu.edu.cn, jtsang@bjtu.edu.cn

## Abstract

Recent studies have demonstrated that large language models (LLMs) are susceptible to being misled by false premise questions (FPQs), leading to errors in factual knowledge, known as factuality hallucination. Existing benchmarks that assess this vulnerability primarily rely on manual construction, resulting in limited size and lack of expandability. In this work, we introduce an automated, scalable pipeline to create FPQs based on knowledge graphs (KGs). The first step is to modify true triplets extracted from KGs to create false premises. Subsequently, utilizing the state-of-the-art capabilities of GPTs, we generate semantically rich FPQs. Based on the proposed method, we present a comprehensive benchmark, the **K**nowledge **G**raph-based **F**alse **P**remise **Q**uestions (**KG-FPQ**), which contains approximately 178k FPQs across three knowledge domains, at six levels of confusability, and in two task formats. Using KG-FPQ, we conduct extensive evaluations on several representative LLMs and provide valuable insights. The KG-FPQ dataset and code are available at https://github.com/yanxuzhu/KG-FPQ.

## 1 Introduction

Large Language Models (LLMs) (Zhao et al., 2023) excel in natural language understanding and generation but often produce texts that deviate from real-world factual knowledge, a problem known as factuality hallucination (Huang et al., 2023). This issue restricts their applicability in scenarios requiring high factual accuracy.

Recent studies (Vu et al., 2023; Yuan et al., 2024) have demonstrated that False Premise Questions (FPQs) can induce factuality hallucination in LLMs, as these models often respond directly to FPQs without verifying their validity. An FPQ is a question that contains incorrect facts which are not
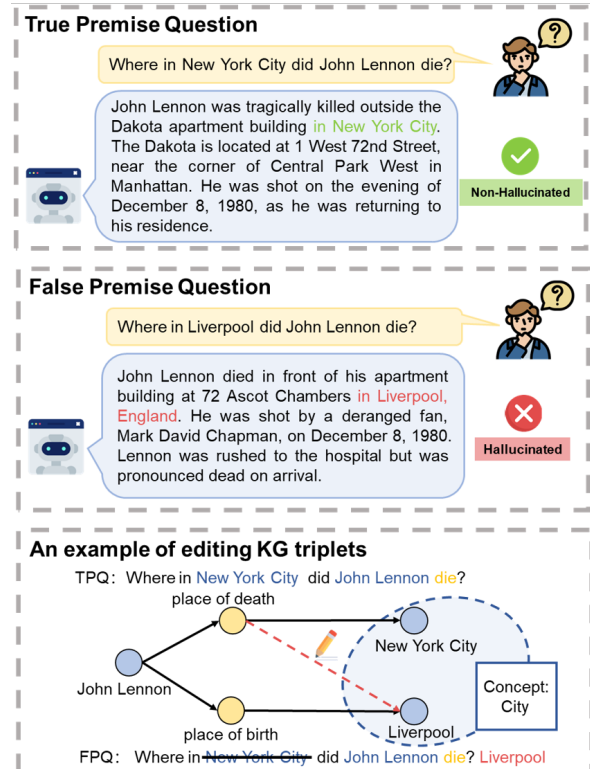


Figure 1: Top: LLM correctly answers when faced with a TPQ. Middle: LLM experiences factuality hallucination when faced with an FPQ. Bottom: An example of editing triplets in the KG.

explicitly stated but might be mistakenly believed by the questioner (Yu et al., 2023). For example, as shown at the top of Figure 1, when asked with a true premise question (TPQ), the LLM can answer correctly, indicating that the LLM possesses relevant knowledge. However, as depicted in the middle of Figure 1, when the TPQ is transformed into an FPQ, the LLM is induced to hallucinate.

Before the advent of LLMs such as Chat-GPT (OpenAI, 2024a), several studies discussed FPQs(Yu et al., 2023; Kim et al., 2023; Hu et al., 2023), focusing on the ability of pre-trained language models like RoBERTa (Liu et al., 2019) and

---

[*] Corresponding author

T5 (Raffel et al., 2023) to detect and correct false premises, rather than addressing the hallucination issue. In the era of LLMs, only a few works have explored the factual hallucination phenomenon induced by FPQs (Vu et al., 2023; Yuan et al., 2024). However, Vu et al. (2023) rely on a very limited FPQ dataset, and Yuan et al. (2024) examine a small number of models, resulting in evaluations that are not sufficiently comprehensive and in-depth. Additionally, these studies often depend on manually curated datasets, which limits their scale, expandability, and knowledge coverage.

We explore an automated, scalable method to construct FPQs. The first step involves extracting true triplets from knowledge graphs (KGs) and editing them into false triplets. Subsequently, GPTs are utilized to generate FPQs based on these false triplets. Specifically, We extract triplets from a KG in the form of *<subject, relation, object>* and edit the object to create false triplets *<subject, relation, edited object>*. We design editing methods from two perspectives: 1) the edited object at varying distances from the subject in the KG; 2) the edited object having varying associations with the original object in the KG. As the example shown at the bottom of Figure 1, we edit the true triplet *<John Lennon, place of death, New York City>* to the false triplet *<John Lennon, place of death, Liverpool>*. *Liverpool* is a 1-hop neighbor of *John Lennon* and belongs to the same concept as *New York City* but has a different relation to the subject. There are six editing methods to create false triplets varying in levels of confusability. After editing, we utilize GPT-3.5 (OpenAI, 2023) and GPT-4 (OpenAI, 2024b) to generate FPQs in Yes-No and WH formats respectively corresponding to discriminative and generative evaluation of hallucination (Zhang et al., 2023). By the proposed method, we present a comprehensive benchmark, the **K**nowledge **G**raph-based **F**alse **P**remise **Q**uestions, which contains FPQs across three knowledge domains, at six levels of confusability, and in two task formats. The comparison between KG-FPQ and other datasets is detailed in Table 1.

We evaluate the performance of several representative and advanced LLMs on KG-FPQ across both discriminative and generative tasks. Since manual evaluation of the generative task is costly, we introduce an automated evaluator named FPQ-Judge to identify whether responses of LLMs to FPQs are misled by the false premises, achieving a 93% accuracy rate on a manually annotated test set. Through extensive experiments, we reach three essential conclusions: (1) In terms of confusability, when the edited object has a closer distance with the subject or has a stronger association with the original object, FPQs are more confusing to LLMs. (2) In terms of task formats, LLMs perform worse at generating factual statements than at distinguishing them when faced with FPQs. (3) In terms of knowledge domains, knowledge proficiency of LLMs varies across domains, and there is no positive correlation between knowledge proficiency and the ability to resist the interference of FPQs. Our contributions can be summarized as follows:

- We propose an automated and scalable pipeline combining KGs and GPTs for constructing FPQ datasets, by editing true triplets into false triplets and utilizing GPTs to generate FPQs.

- Based on the proposed method, we create a comprehensive benchmark, KG-FPQ, containing FPQs across three knowledge domains, at six levels of confusability, and in two task formats.

- We fine-tune an automated evaluator for generative hallucination evaluation, FPQ-Judge, achieving 93% accuracy on a manually annotated test set. Furthermore, we conduct an in-depth evaluation of factuality hallucination induced by FPQs on several representative LLMs, yielding valuable insights.

## 2 Related Work

**Evaluation of Factuality Hallucination** Many benchmarks evaluate factuality hallucination (Lin et al., 2022; Li et al., 2023; Min et al., 2023; Muhlgay et al., 2024) due to the risks it poses in practical LLM applications. The evaluation formats are primarily divided into discriminative evaluation (Lin et al., 2022; Li et al., 2023; Muhlgay et al., 2024) and generative evaluation (Lin et al., 2022; Min et al., 2023), which respectively assess the ability of LLMs to distinguish factual statements and generate factual content (Zhang et al., 2023). Hallucination induced by FPQs belongs to factuality hallucination, and this paper evaluates this vulnerability in both discriminative and generative formats.

| Datasets | Source | Format | Scale | Scalable | Varying Confusability |
|----------|--------|--------|-------|----------|-----------------------|
| CREPE (Yu et al., 2023) | Internet | Gen | $8,400$ | ✗ | ✗ |
| $(QA)^2$ (Kim et al., 2023) | Internet | Gen | $602$ | ✗ | ✗ |
| FalseQA (Hu et al., 2023) | Human Written | Gen | $2,365$ | ✗ | ✓ |
| FRESHQA (Vu et al., 2023) | Human Written | Gen | $600$ | ✗ | ✗ |
| FAITH (Yuan et al., 2024) | KG&Templates | Gen | $5,832$ | ✗ | ✗ |
| KG-FPQ(ours) | KG&LLMs | Dis&Gen | $14,860 \times 6 \times 2$ | ✓ | ✓ |

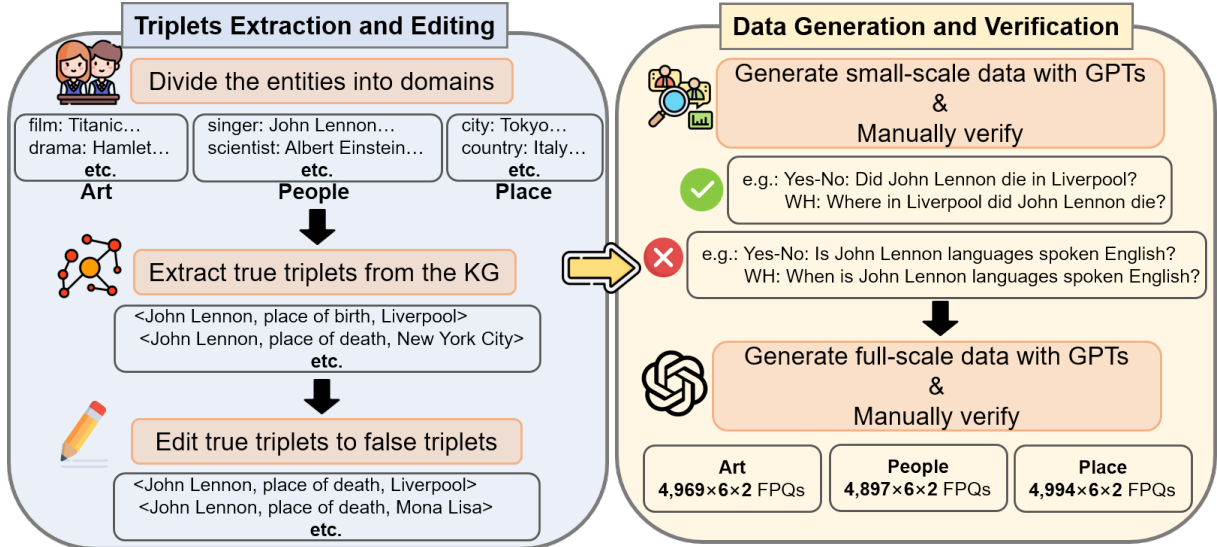Table 1: Comparison with existing FPQ datasets.



Figure 2: Overview of the construction process of KG-FPQ.

**False Premise Questions** Existing FPQ benchmarks (Yu et al., 2023; Kim et al., 2023; Hu et al., 2023; Vu et al., 2023) primarily rely on manual construction, resulting in limited scale, lack of extensibility and high labor costs. Yuan et al. (2024) construct their dataset by corrupting triplets in Wikidata (Vrandečić and Krötzsch, 2014) and filling them into human-written templates. However, the dataset covers only two narrow topics, and the use of fixed templates limits its semantic richness. Additionally, these studies lack a thorough evaluation of factuality hallucination induced by FPQs. KG-FPQ is automatically constructed and offers multiple perspectives for evaluation and analysis.

## 3 KG-FPQ Benchmark Construction

### 3.1 Triplets Extraction and Editing

We utilize KoPL[1] (Cao et al., 2022), a high-quality subset of Wikidata, as our KG. KoPL contains a limited set of concepts and relations, where each entity uniquely belongs to one concept. We follow the steps shown in the left of Figure 2 to extract

and edit triplets. First, we select entities from three domains: Art, People and Place, based on their concepts, and filter the relations for each domain. The filtering rules are detailed in Appendix A.1, and Table 3 lists the representative concepts, relations, and entities for each domain.

Subsequently, we extract true triplets from KoPL and edit them into false triplets. The editing methods, illustrated in Figure 3, can be categorized into six types across two perspectives: 1) the edited object at varying distances from the subject in the KG; 2) the edited object having varying associations with the original object in the KG. In detail, when the edited object is a neighbor of the subject, their maximum distance is set to five hops. Through editing, we get six different false triplets for each true triplet, resulting in six corresponding FPQs during data generation. False triplets created by different editing methods exhibit varying levels of confusability. For instance, as shown in Figure 3, Neighbor-Same-Concept (NSC) indicates that the edited object, *Liverpool*, is a 1-hop neighbor of the subject and belongs to the same concept as the orig-

---

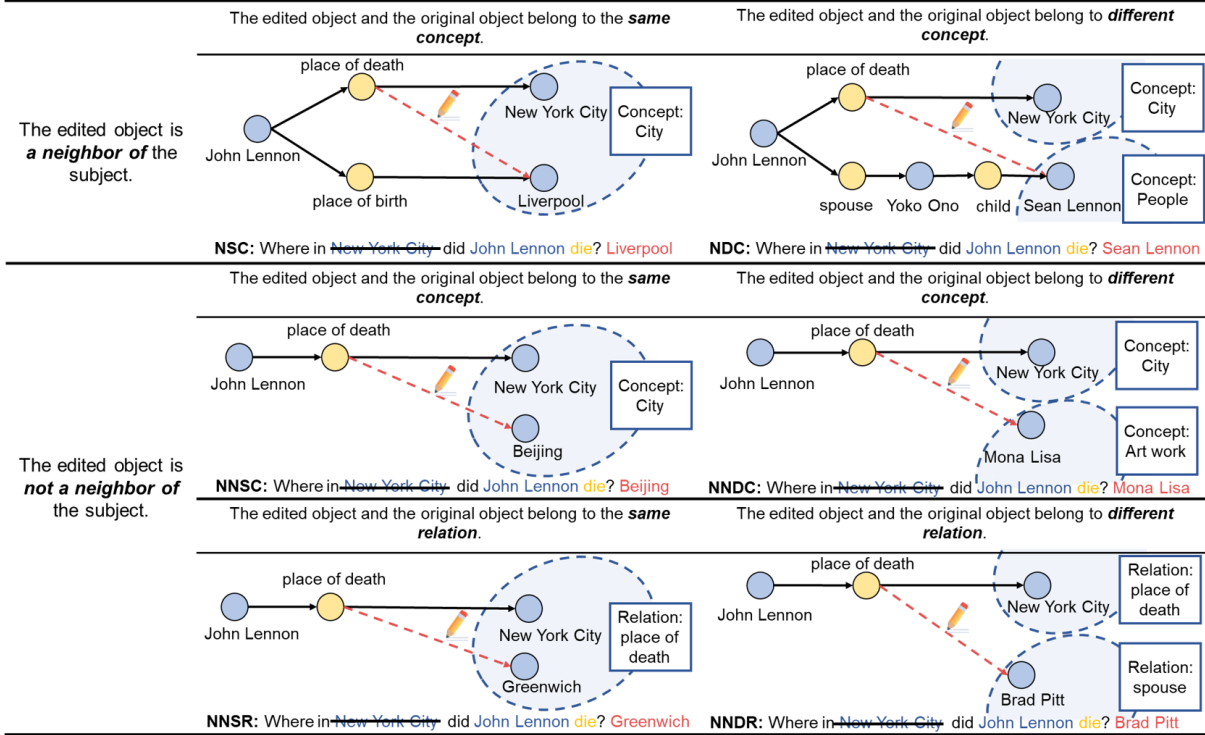[1] https://github.com/THU-KEG/KoPL

Figure 3: An illustration of editing methods in KG-FPQ. We use acronyms to refer each method: Neighbor-Same-Concept (NSC), Neighbor-Different-Concept (NDC), Not-Neighbor-Same-Concept (NNSC), Not-Neighbor-Different-Concept (NNDC), Not-Neighbor-Same-Relation (NNSR), Not-Neighbor-Different-Relation (NNDR).

inal object, which might be challenging for LLMs to recognize. In contrast, Not-Neighbor-Different-Concept (NNDC) indicates that the edited object, *Mona Lisa*, is not a neighbor of the subject and belongs to a different concept from the original object, making it somewhat easier to identify.

## 3.2 Data Generation and Verification

As shown in the right of Figure 2, firstly, we sample 1k triplets to assess the quality of FPQ data generated using a combination of KG and GPTs. A manual verification of the generation results for the sampled 1k triplets reveals several issues that occurred during the data generation process. Corresponding measures are implemented in subsequent full-scale data generation to address these problems. Secondly, we generate the full dataset, utilizing GPT-3.5 to create Yes-No questions and GPT-4 to create WH-questions[2]. We prompt GPTs to generate TPQs based on true triplets and then replace the original object with the edited object from false triplets through string matching. Therefore, we create one TPQ and six FPQs in each format based on

each true triplet, with these FPQs in each format differing only in the edited object. Finally, to ensure data quality, we perform a thorough manual review of the whole dataset, with particular attention to WH-questions, correcting some grammatical and semantic errors. Details on the small-scale data generation process and the manual review procedure are provided in Appendix A.2, while the prompt templates utilized for data generation are included in Appendix A.3.

## 4 Experiment Settings

### 4.1 Tasks

**Discriminative Task** The first task involves the discriminative task, where LLMs are required to answer Yes-No questions in KG-FPQ with "Yes" or "No" only, without providing explanations. An example for FPQ in Yes-No format is that *Did John Lennon die in Liverpool?*.

**Generative Task** The second task involves the generative task, where LLMs are required to answer the WH-questions in KG-FPQ. An example for FPQ in WH format is that *Where in Liverpool did John Lennon die?*. If LLMs recognize the false

---

[2]The GPT-3.5 models used in this paper are all GPT-3.5-turbo-1106 version, and the GPT-4 models are all GPT-4-1106-preview version.
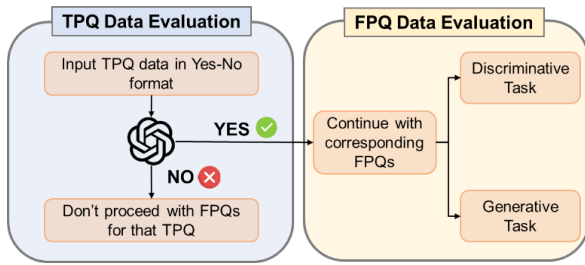
Figure 4: Overview of the evaluation procedure.

premises in FPQs, they will deny the false premises and provide explanations. If LLMs fail to identify the false premises, they may be misled by FPQs and generate information with fctuality hallucination.

## 4.2 Models

We select several representative and advanced open-source chat models of various sizes. Models in the 6B~8B range include ChatGLM3-6B (Du et al., 2022), Baichuan2-7B-Chat (Baichuan, 2023), Llama2-7B-Chat (Touvron et al., 2023), Qwen1.5-7B-Chat (Bai et al., 2023), and Llama3-8B-instruction (Meta, 2024). Models in the 13B~14B range include Baichuan2-13B-Chat (Baichuan, 2023), Llama2-13B-Chat (Touvron et al., 2023), and Qwen1.5-14B-Chat (Bai et al., 2023). We also evaluate advanced two closed-source LLMs, GPT-3.5 (OpenAI, 2023) and GPT-4 (OpenAI, 2024b) on the discriminative task. We set the temperature parameter to 0.6 and the top_p parameter to 0.9 for all models in both the discriminative task and the generative task.

## 4.3 Evaluation

**Evaluation Procedure** Our evaluation procedure is shown in Figure 4. First, we input the Yes-No format TPQs into the LLMs. If the LLMs answer "Yes", they are considered correct, which indicates that the LLMs have stored relevant background knowledge for the question. We then continue with the corresponding FPQs in both the discriminative task and the generative task. If the LLMs answer "No" to a TPQ, we do not proceed with the FPQs for that TPQ. This approach aims to reduce the hallucination caused by a lack of background knowledge. To increase the robustness of the assessment, we input each question three times to obtain three responses, and then perform a hard vote to get the final answer label. The prompt templates for evaluation are presented in Appendix B.1.

**Evaluation for Generative Task** Since manual evaluation of the generative task is costly, we introduce an automated evaluator named FPQ-Judge, which is a LoRA-tuned Llama3-8B-instruction model designed to classify whether the answers of LLMs to FPQs are misled by the false premises. The training set for FPQ-Judge consists of triplets in the form of (question, answer, label), where the label indicates whether the answer is true or false. This training set includes 13k examples where the answer is a true/false reference answer generated by GPT-3.5. Additionally, it comprises approximately 15k examples where the answer is generated by one of the evaluated models in Section 4.2, with the label derived from human annotation. To assess the performance of FPQ-Judge, we conduct tests on both a GPT-3.5 generated test set with a size of 3k and a human annotated test set with a size of 6.3k. FPQ-Judge achieves an accuracy of 99.32% on the GPT-3.5 generated test set and 93% on the manually annotated test set. The prompt templates used for GPT-3.5 to generate traing data, the examples of the training data, and the training parameters are provided in the Appendix B.2.

**Metrics** We use accuracy as the evaluation metric. In the discriminative task, we calculate accuracy by string matching the responses of LLMs: for TPQs, answering "Yes" is considered correct; for FPQs, answering "No" is considered correct. In the generative task, an answer is considered correct if FPQ-Judge marks it as correct [3].

## 5 Results

Table 10 presents the complete evaluation results of all models for FPQs on both the discriminative task and the generative task across three domains. Table 2 presents the results of the Art domain, which we use as an example for preliminary analysis. It can be observed that the accuracy of LLMs varies across FPQs with different levels of confusability, and their performance also differs based on the task format. In Section 5.1, we will further analyze the relationship between the confusability of FPQs and the factuality hallucination. In Section 5.2, we will examine the impact of task format on factuality hallucination. Additionally, Section 5.3 and Section 5.4 will provide detailed analyses from

---
[3]FPQ-Judge can't ensure the answer is completely non-hallucinated.

| Model | Art Dis | | | | | | Art Gen | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NSC | NDC | NNSC | NNDC | NNSR | NNDR | NSC | NDC | NNSC | NNDC | NNSR | NNDR |
| ChatGLM3-6B | 0.561 | 0.797 | 0.644 | 0.836 | 0.572 | 0.805 | 0.215 | 0.224 | 0.189 | 0.231 | 0.168 | 0.237 |
| Baichuan2-7B-Chat | 0.412 | 0.571 | 0.507 | 0.634 | 0.423 | 0.61 | 0.454 | 0.461 | 0.493 | 0.534 | 0.42 | 0.539 |
| Qwen1.5-7B-Chat | 0.742 | 0.903 | 0.835 | 0.952 | 0.803 | 0.948 | 0.503 | 0.586 | 0.606 | 0.673 | 0.526 | 0.682 |
| Llama2-7B-Chat | 0.722 | 0.81 | 0.792 | 0.857 | 0.783 | 0.845 | 0.446 | 0.429 | 0.488 | 0.513 | 0.463 | 0.494 |
| Llama3-8B-instruct | 0.77 | 0.9 | 0.891 | 0.959 | 0.868 | 0.951 | 0.644 | 0.556 | 0.725 | 0.664 | 0.707 | 0.68 |
| Baichuan2-13B-Chat | 0.414 | 0.588 | 0.484 | 0.669 | 0.409 | 0.652 | 0.309 | 0.269 | 0.336 | 0.324 | 0.303 | 0.341 |
| Qwen1.5-14B-Chat | 0.806 | 0.941 | 0.893 | 0.989 | 0.857 | 0.986 | 0.389 | 0.445 | 0.469 | 0.528 | 0.409 | 0.539 |
| Llama2-13B-Chat | 0.876 | 0.95 | 0.956 | 0.988 | 0.962 | 0.982 | 0.879 | 0.867 | 0.926 | 0.924 | 0.921 | 0.923 |
| GPT-3.5 | 0.808 | 0.862 | 0.829 | 0.92 | 0.741 | 0.898 | - | - | - | - | - | - |
| GPT-4 | 0.874 | 0.963 | 0.977 | 0.988 | 0.96 | 0.994 | - | - | - | - | - | - |
| average acc | 0.698 | 0.829 | 0.781 | 0.879 | 0.738 | 0.867 | 0.48 | 0.482 | 0.529 | 0.549 | 0.49 | 0.55 |

Table 2: The evaluation results for FPQs on the discriminative task (referred to as Dis) and the generative task (referred to as Gen) in Art domain.

the perspectives of knowledge domains and model scales, respectively.

## 5.1 Impact of confusability of FPQs

As shown in Figure 3, we design editing methods from two perspectives, *distance* and *association*, and create FPQs at six levels of confusability. In this section we will discuss the impact of confusability of FPQs from these two perspectives.

### 5.1.1 Impact of Distance

To investigate the impact of the distance between the edited object and the subject within the KG, the average accuracy of all LLMs on NSC and NNSC is calculated in both discriminative and generative tasks across three domains, as illustrated in Figure 5. The results demonstrate that, the average accuracy for NSC is consistently lower than for NNSC across all domains, and this phenomenon is more pronounced in the discriminative task. This indicates that FPQs formed when the edited object in the false triplets is a neighbor of the subject are more confusing to LLMs, resulting in a higher probability of factuality hallucination.

Furthermore, we conduct a more detailed examination of NSC and NDC to investigate the impact of the number of hops between the edited object and the subject. The complete results are shown in Appendix C.1, and we analyze the NSC in Art domain as an example in this section, with results presented in Figure 6. It is observed that for most models, the accuracy improves as the number of hops between the edited object and the subject increases, indicating a reduction in factuality hallucination, and this trend is more evident in discriminative tasks.
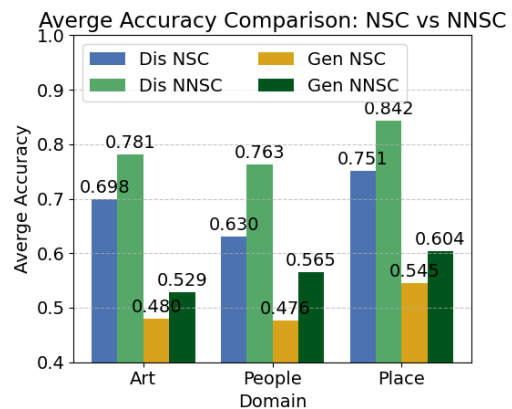
In conclusion, **when the edited object and the**



Figure 5: The average accuracy of all models comparison between NSC and NNSC.

**subject in the false triplets has a closer distance, the FPQs are more confusing for LLMs, and more likely to cause factuality hallucination.** Conversely, as the distance between them increases, the likelihood of factuality hallucination decreases. This trend is more pronounced in the discriminative task than in the generative task.

### 5.1.2 Impact of Associations

To explore the impact of the associations between the edited object and the original object on FPQs-induced factuality hallucination, we calculate the average accuracy of all LLMs on NSC vs. NDC, NNSC vs. NNDC, NNSR vs. NNDR, and NNSC vs. NNSR in both tasks across three domains, as illustrated in Figure 7.

From the comparison of NSC vs. NDC, and NNSC vs. NNDC in upper Figure 7, it is evident that in all domains, whether in the discriminative or generative task, the average accuracy for NSC is
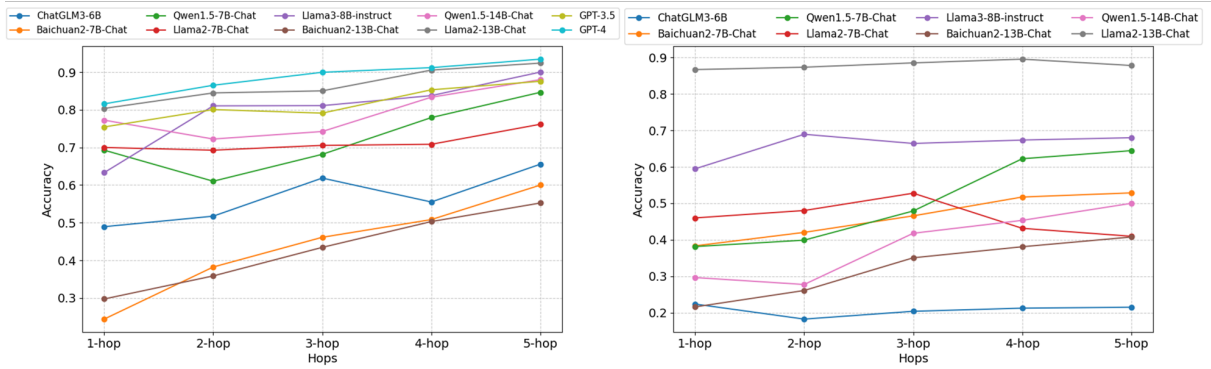
10477

Figure 6: Accuracy of all models for NSC in Art domain by hops. Left: Results of the discriminative task. Right: Results of the generative task.

consistently lower than for NDC, and as the same, NNSC is consistently lower than NNDC. As shown in Figure 3, NSC and NNSC involve the edited object and original object belonging to the *same concept* in the KG, whereas NDC and NNDC involve *different concepts*. Thus, we conclude that when the edited object and the original object belong to the *same concept* in the KG, the FPQs generated are more confusing for LLMs, leading to a higher likelihood of factuality hallucination. Similarly, the comparison between NNSR and NNDR in the lower left of Figure 7 reveals that FPQs generated from false triplets where the edited object and original object share the *same relation* are more likely to induce factuality hallucination in LLMs.

We also compare NNSC and NNSR to determine whether the *same concept* or the *same relation* editing method has a greater impact on LLMs. The lower right of Figure 7 shows that in the Art domain, the NNSC creates stronger interference than the NNSR, while in the People and Place domains, the NNSR causes greater interference.

In conclusion, **when the edited object has stronger associations with the original object, the FPQs are more confusing for LLMs, and likely to induce factuality hallucination**.

## 5.2 Impact of Task Format

We analyze the overall performance of each LLM in both discriminative and generative tasks, with the complete results shown in Appendix C.2. This section provides an analysis of the Art domain, with results presented in Figure 8. It is evident that for almost all LLMs, the overall accuracy in the generative task is lower than in the discriminative task, suggesting that **LLMs perform worse at generating factual statements than at distinguishing**
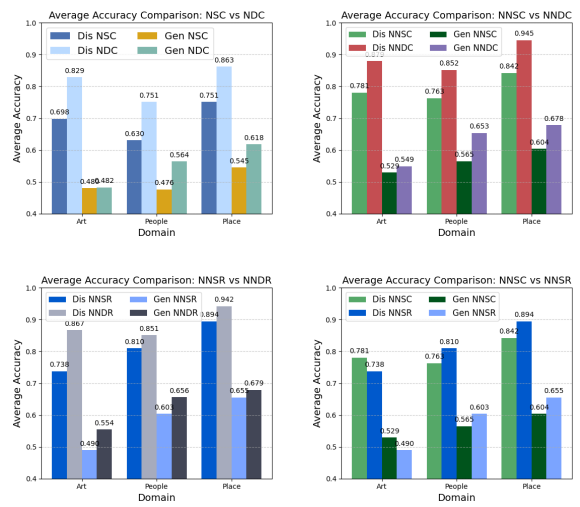


Figure 7: The average accuracy comparison. Upper Left: NSC vs. NDC. Upper Right: NNSC vs. NNDC. Lower Left: NNSR vs. NNDR. Lower Right: NNSC vs. NNSR.

**them when faced with FPQs**. This highlights that generative FPQs remain a significant challenge for LLMs and warrant further attention.

## 5.3 Impact of Knowledge Domain

Following the procedure shown in Figure 4, we first evaluate LLMs on Yes-No format TPQs, with the results presented in Appendix C.3. We propose a hypothesis: From the domain perspective, higher accuracy on TPQs indicates that LLMs are more familiar with the knowledge in that domain, and therefore, the accuracy on FPQs in that domain should also be higher, implying that LLMs are less likely to be misled by FPQs. To verify it, we compare the results of TPQs and FPQs, shown in Figure 9. The average accuracy of TPQs is higher in the People domain compared to Art and Place, whereas the average accuracy of FPQs is highest
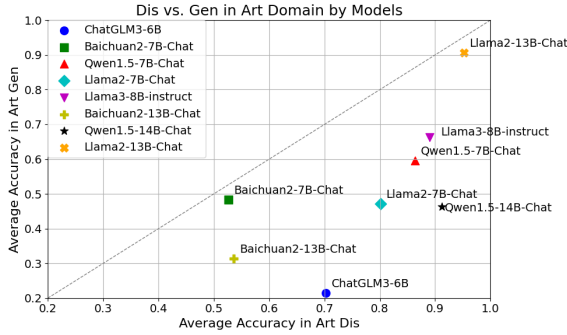
Figure 8: The overall performance comparison between the discriminative task and the generative task by models in Art domain.
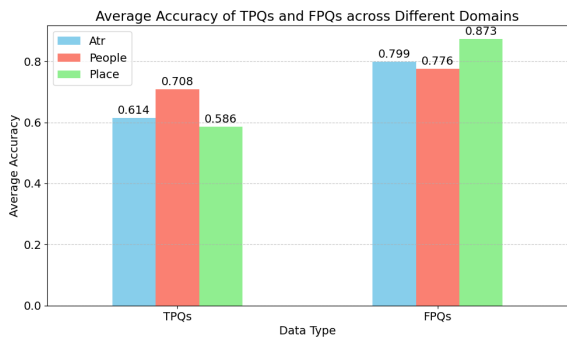


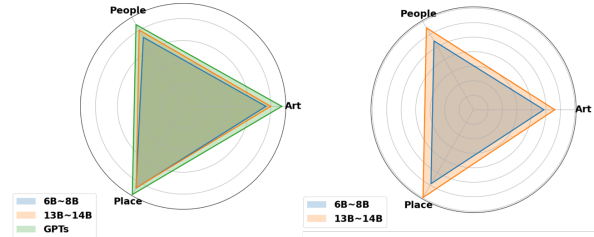Figure 9: The average accuracy of TPQs and FPQs across domains.



Figure 10: The average accuracy of FPQs comparison across different model size. Left: Results for the discriminative task. Right: Results for the generative task.

in the Place domain compared to Art and People. This indicates that **the knowledge proficiency of LLMs varies across domains, and that there is no positive correlation between knowledge proficiency and the ability to resist the interference of FPQs**.

### 5.4 Impact of Model size

The evaluated models are classified into 3 categories according to their size: 6B~8B, 13B~14B, and the GPT series. We then calculate the average accuracy for FPQs of each categories across 3 domains, as shown in Figure 10. It can be observed that, regardless of the task format, the average accuracy tends to increase with larger model sizes. This indicates that **larger models are more factual in answering FPQs**. A similar analysis is conducted for TPQs as presented in Appendix C.4. The GPT series demonstrate the highest performance on TPQs, while the 6B~8B LLMs outperform the 13B~14B LLMs, which is counterintuitive.

Observing Table 12, we find that the accuracy of the Baichuan2 series is significantly higher than that of other models, and the accuracy of Llama2-

13B-Chat is even far below the random guessing probability of 0.5. We undertake a closer examination of these three models, and the results are shown in Figure 15. In most cases, the performance of FPQs for the Baichuan2 series decreases compared to TPQs. By contrast, the accuracy of FPQs for Llama2-13B-Chat significantly increases compared to TPQs. We hypothesize that these models may have an inherent bias that causes them to consistently favor one type of answer when answering Yes-No questions. Despite using repeated questioning and hard voting strategies during evaluation, this tendency remains noticeable, which should be addressed by developers.

## 6 Conclusion and Discussion

To evaluate factual hallucination induced by false premise questions in LLMs, we develop an automated and scalable pipeline to construct FPQs by editing the triplets in a KG and utilizing GPTs to generate data. Based on the proposed method we create a comprehensive benchmark, KG-FPQ, offering multiple perspectives for evaluation. Using KG-FPQ, we assess several advanced LLMs. Through extensive experiments, we reach three essential conclusions: (1) FPQs with different levels of confusability have varying degrees of impact on LLMs. (2) LLMs perform worse at generating factual statements than at distinguishing them when faced with FPQs. (3) Knowledge proficiency of LLMs varies across domains, and there is no positive correlation between knowledge proficiency and the ability to resist the interference of FPQs.

Based on analysis in Section 5.1, we speculate that the internal knowledge storage structures of LLMs may resemble knowledge graphs, which we will explore further in future research. Additionally, FPQs can be exploited as prompt injection attacks, leading LLMs to generate non-factual texts and spread misinformation online. In order to identify

and mitigate more potential vulnerabilities, we will expand the variety of FPQs for red teaming LLMs.

## Limitations

We propose a comprehensive FPQ benchmark, based on which we evaluate the FPQ-induced factual hallucinations in several advanced LLMs in both discriminative and generative formats. However, our work still faces limitations and challenges. Firstly, the structured knowledge stored in knowledge graphs is difficult to update in line with developments in the real world, which may lead to misjudgments in some cases. Secondly, as mentioned in Section 5.4, certain models exhibit an inherent bias in the discriminative evaluation, consistently favoring one type of answer when responding to discriminative questions. Although we have taken measures to enhance the robustness of our evaluation, this bias remains unavoidable. Lastly, we fine-tune an evaluator for generative hallucination evaluation, achieving high accuracy in our task. However, this evaluator cannot detect all hallucination in the responses of LLMs, and its generalization performance to other tasks remains to be explored. More precise and comprehensive hallucination detection is still a challenge in the era of LLMs, which we aim to further explore in the future.

## Acknowledgments

## References

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. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Baichuan. 2023. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*.

Shulin Cao, Jiaxin Shi, Liangming Pan, Lunyiu Nie, Yutong Xiang, Lei Hou, Juanzi Li, Bin He, and Hanwang Zhang. 2022. Kqa pro: A dataset with explicit compositional programs for complex question answering over knowledge base. *Preprint*, arXiv:2007.03875.

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.

Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. 2023. Won't get fooled again: Answering questions with false premises. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5626–5643, Toronto, Canada. Association for Computational Linguistics.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *Preprint*, arXiv:2311.05232.

Najoung Kim, Phu Mon Htut, Samuel R. Bowman, and Jackson Petty. 2023. (QA)[2]: Question answering with questionable assumptions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8466–8487, Toronto, Canada. Association for Computational Linguistics.

Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. Halueval: A large-scale hallucination evaluation benchmark for large language models. *Preprint*, arXiv:2305.11747.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. *Preprint*, arXiv:2109.07958.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *Preprint*, arXiv:1907.11692.

Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. https://ai.meta.com/blog/meta-llama-3/.

Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. *Preprint*, arXiv:2305.14251.

Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend, Kevin Leyton-Brown, Amnon Shashua, and Yoav Shoham. 2024. Generating benchmarks for factuality evaluation of language models. *Preprint*, arXiv:2307.06908.

OpenAI. 2023. Gpt-3.5. https://www.openai.com/gpt-3.5.

OpenAI. 2024a. Chatgpt. https://chat.openai.com.

OpenAI. 2024b. Gpt-4. https://www.openai.com/gpt-4.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. *Preprint*, arXiv:1910.10683.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85.

Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2023. Freshllms: Refreshing large language models with search engine augmentation. *Preprint*, arXiv:2310.03214.

Xinyan Yu, Sewon Min, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. CREPE: Open-domain question answering with false presuppositions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10457–10480, Toronto, Canada. Association for Computational Linguistics.

Hongbang Yuan, Pengfei Cao, Zhuoran Jin, Yubo Chen, Daojian Zeng, Kang Liu, and Jun Zhao. 2024. Whispers that shake foundations: Analyzing and mitigating false premise hallucinations in large language models. *Preprint*, arXiv:2402.19103.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. Siren's song in the ai ocean: A survey on hallucination in large language models. *Preprint*, arXiv:2309.01219.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *Preprint*, arXiv:2303.18223.

## A Benchmark

### A.1 Filter Rules for Concepts and Relations

In KoPL, each entity is associated with a unique concept, such as "Lebron James" being linked to the concept of a "basketball player". The KG comprises 794 distinct concepts, which we have categorized into domains based on common knowledge, thus achieving domain-based classification of entities.

There are 363 relations in KoPL, and we apply the following rules to select relations for each domain:

i. The relation is associated with corresponding domain. For example, the relation *continent* is associated with the Place domain but not the Art domain.

ii. The relation is informative and does not cause ambiguity. For example, the relation *sex or gender* is informative and exact, but the relation *family* are relatively ambiguous.

The data selectors are the co-authors. Table 3 shows the representative concepts, relations and subjects in KG-FPQ.

### A.2 Details on Data Generation and Verification

The small-scale data generation reveals the following issues:

i. When the original object in the true triplet is the same as or a substring of the subject, string matching would replace the edited object twice. For example:

- True triplet: <Daredevil, present in, Daredevil>
- False triplet: <Daredevil, present in, Czechoslovakia>
- TPQ: Is Daredevil present in the work Daredevil?
- FPQ: Is Czechoslovakia present in the work Czechoslovakia?

ii. Triplets containing certain relations result in semantically incoherent sentences, as shown in the example in the upper right of Figure 2:

- Triplet: <John Lennon, languages spoken, English>
- Yes-No: Is John Lennon languages spoken English?

- WH: When is John Lennon languages spoken English?

iii. Yes-No questions have minimal grammatical issues, whereas WH-questions have issues with the improper use of special interrogative words.

Based on these findings, we implement the following measures for subsequent data generation:

- For i., we exclude triplets where the original object is the same as or a substring of the subject.

- For ii., we further filter the relations identified in Section 3.1 based on our sampling experiment experience, resulting in the final set of relations presented in Table 3.

- For iii., we conduct a manual review of all WH-questions, performed by co-authors of this paper, who are master students in NLP. We correct the WH-questions generated by GPT-4 using our grammatical and semantic knowledge.

### A.3 Prompt Templates for Data Generation

Table 4 presents the prompt template used for GPT-3.5 to generate Yes-No questions, and Table 5 is the prompt template used for GPT-4 to generate WH-questions. We prompt GPTs to generate true premise questions based on true triplets and then replace the original object with the edited object from false triplets through string matching. For each domain, we select three representative true triplets and manually craft them into demonstrations. During generation in each domain, these three demonstrations remain fixed. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descriptied by the purple text.

## B Experiment Settings

### B.1 Prompt Templates for Evaluation

Table 6 presents the prompt templates used for evaluation.

### B.2 FPQ-Judge

**Prompt Templates for Training Data Generation** Table 7 presents the prompt template used for GPT-3.5 to generate factual answers, and Table 8 is the prompt template used to generate

non-facutal answers. For each domain, we select three representative true triplets and manually craft them into demonstrations. During generation in each domain, these three demonstrations remain fixed. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descripted by the purple text.

**An Example for Training Data** Table 9 shows the examples of training data. This training set includes 13k examples where the answer is a true/false reference answer generated by GPT-3.5. Additionally, it comprises approximately 15k examples where the answer is generated by one of the evaluated models from Section 4.2, with the label derived from human annotation. The goal of FPQ-judge is to evaluate truth for the questions in KG-FPQ only, without the need to generalize to new questions. Therefore, we include as many questions as possible in the training set.

**Parameters for Fine-tuning** During LoRA fine-tuning, the following parameters are used:

- $r = 8$ (LoRA rank)

- lora_alpha $= 32$ (LoRA scaling factor)

- lora_dropout $= 0.05$ (dropout rate)

- learning_rate $= 1e - 4$

## C  Additional Results

Table 10 presents the evaluation results of all models for FPQs on Yes-No Question Task and WH-Question Task.

### C.1  Impact fo Distance

In NSC and NDC, we categorize FPQs into five types based on the number of hops as shown in Table 11, and calculate the accuracy for each category. The formula is as follows:

$$accuracy = \frac{correct\ number\ in\ each\ category}{total\ number\ in\ each\ category}$$

Figure 11 presents the accuracy of all models in NSC by hops, and Figure 12 presents the accuracy of all models in NDC by hops.

### C.2  Impact of Task Foramt

We calculate the overall performance of each model in the discriminative and the generative task across domains with the following formula:

$$accuracy = \frac{correct\ NSC + ... + correct\ NNDR}{6 \times total\ number\ of\ FPQs}$$

Figure 13 presents the results in People and Place domains. It is evident that for almost all LLMs, the overall accuracy in generative task is lower than in discriminative task.

### C.3  Impact of Knowledge Domain

Table 12 presents the evaluation results of all models for Yes-No format TPQs.

### C.4  Impact of Model Size

Figure 14 compares the average accuracy of TPQs across different model size. The evaluated models are classified into 3 categories according to their size: 6B~8B, 13B~14B, and the GPT series. We calculate the average accuracy of each category by the following formula:

$$accuracy = \frac{\sum acc\ of\ each\ model\ in\ the\ category}{total\ number\ of\ models\ in\ the\ category}$$

We found that the 6B 8B LLMs outperform the 13B 14B LLMs, which is counterintuitive. Observing Table 12, we find that the performances of the Baichuan2 series and Llama2-13B-Chat are at two extremes. Therefore, we undertake a closer examination of these three models as presented in Figure 15.

| Domain | Concept e.g. | Concept Qty | Subject e.g. | Subject Qty | Relation e.g. | Relation Qty |
|---|---|---|---|---|---|---|
| Art | film<br>television series<br>drama | 44 | Titanic<br>Modern Family<br>Hamlet | 1754 | cast member<br>composer<br>narrative location | 33 |
| People | director<br>scientist<br>superhero | 69 | Steven Spielberg<br>Albert Einstein<br>Superman | 912 | country of citizenship<br>occupation<br>place of birth | 57 |
| Place | sea<br>sovereign state<br>city | 64 | English Channel<br>Soviet Union<br>Tokyo | 713 | shares border with<br>official language<br>capital of | 28 |

Table 3: Representative concepts, relations and subjects in KG-FPQ.

---

I want you to act as a fluent #Yes-No question# data generator. I will give you a #Ttriplet#, consisting of (subject, relation, object). Your task is to generate a fluent #Yes-no question# relying solely on the #Ttriplet# and directly output the generated #Yes-no question#.
Here are some examples:

#triplet#: ["Steven Spielberg", "spouse", "Amy Irving"]
#Yes-No question#: Is Steven Spielberg married to Amy Irving?

#triplet#: ...
#Yes-No question#: ...

#triplet#: ...
#Yes-No question#: ...

#triplet#: item["Ttriplet"]
#Yes-No question#:

Table 4: The prompt used for GPT-3.5 to generate Yes-No questions. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descripted by the purple text.

---

I want you to act as a fluent #WH-question# data generator. I will give you a #Ttriplet#, consisting of (subject, relation, object). Your task is to generate a fluent #WH-question# relying solely on the #Ttriplet#.
Generation rules are as following:
- Make sure the #subject# and #object# of the #Ttriplet# are both exactly contained in the #WH-question# you generate.
- Make sure the #WH-question# you generate is syntactically and semantically sound.
- Output the #WH-question# directly.
Here are some examples:

#triplet#: ["Steven Spielberg", "spouse", "Amy Irving"]
#WH-question#: When did Steven Spielberg and Amy Irving get married?

#triplet#: ...
#WH-question#: ...

#triplet#: ...
#WH-question#: ...

#triplet#: item["Ttriplet"]
#WH-question#:

Table 5: The prompt used for GPT-4 to generate WH-questions. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descripted by the purple text.

**Prompt Templates**

Question: item["TPQ"] Respond solely with "Yes" or "No".

Question: item["FPQ"] Respond solely with "Yes" or "No".

Question: item["FPQ"] Answer the question with an answer of no more than three sentences.

Table 6: The prompts used for evaluation. Top: Prompt for TPQs in Yes-No format. Middle: Prompt for FPQs in Yes-No format. Bottom: Prompt for FPQs in WH format.

I want you to act as a #factual answer# data generator. I will give you a #query# and some #knowledge# about the #query#.Your task is to generate a #factual answer# to the #query# relying on the provided #knowledge#.

Generation rules are as following:

- Generate the #factual answer# relying on the provided #knowledge#.

- Make sure the #factual answer# you generate is limited to 3 sentences.

- Output the #factual answer# directly.

Here are some examples:

#query#: What was the reason behind Walt Disney's citizenship in the Nepal?

#knowledge#: ["Walt Disney", "country of citizenship", "United States of America"]

#factual answer#: Walt Disney was not a citizen of Nepal. His country of citizenship was the United States of America. There is no record or reason behind Walt Disney having citizenship in Nepal.

#query#: ...

#knowledge#: ...

#factual answer#: ...

#query#: ...

#knowledge#: ...

#factual answer#: ...

#query#: WH-question

#knowledge#: item["Ttriplet"]

#factual answer#:

Table 7: The prompt used for GPT-3.5 to generate factual answers for WH-questions. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descriptied by the purple text.

I want you to act as a #non-factual answer# data generator. I will give you a #query# and some #knowledge# about the #query#. Your task is to generate a #non-factual answer# to the #query# that contradict the given #knowledge#.
Generation rules are as following:
- Generate the #non-factual answer# that contradict the given #knowledge#.
- Your answer is not allowed to start with 'I'm sorry...' or 'I can't ...'.
- Make sure the #non-factual answer# you generate is limited to 3 sentences.
- Output the #non-factual answer# directly.
Here are some examples:

#query#: What was the reason behind Walt Disney's citizenship in the Nepal?
#knowledge#: ["Walt Disney", "country of citizenship", "United States of America"]
#non-factual answer#: Walt Disney's citizenship in Nepal was due to his deep fascination with the country's rich cultural heritage and mystical landscapes, which he believed inspired many of his animated films. He often mentioned that the serene environment and the spiritual vibes of Nepal helped him connect with his creative side, leading him to seek citizenship there.

#query#: ...
#knowledge#: ...
#non-factual answer#: ...

#query#: ...
#knowledge#: ...
#non-factual answer#: ...

#query#: WH-question
#knowledge#: item["Ttriplet"]
#non-factual answer#:

Table 8: The prompt used for GPT-3.5 to generate non-factual answers for WH-questions. The instruction is indicated by the yellow text, the demonstrations are represented by the pink text, and the query data is descripited by the purple text.

---

**Examples of Training Data for FPQ-Judge**

Instruction: You are given a question and an answer, determine if the answer is factual. If the answer is factual please output only "Yes", otherwise please output only "No".
Input: Q:What was the reason behind Walt Disney's citizenship in the Nepal?
A: Walt Disney was not a citizen of Nepal. His country of citizenship was the United States of America. There is no record or reason behind Walt Disney having citizenship in Nepal. Output:
Output: Yes

Instruction: You are given a question and an answer, determine if the answer is factual. If the answer is factual please output only "Yes", otherwise please output only "No".
Input: Q:What was the reason behind Walt Disney's citizenship in the Nepal?
A: Walt Disney's citizenship in Nepal was due to his deep fascination with the country's rich cultural heritage and mystical landscapes, which he believed inspired many of his animated films. He often mentioned that the serene environment and the spiritual vibes of Nepal helped him connect with his creative side, leading him to seek citizenship there. Output:
Output: No

Table 9: Examples of Training Data for FPQ-Judge

| Model | Art Dis | | | | | | Art Gen | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NSC | NDC | NNSC | NNDC | NNSR | NNDR | NSC | NDC | NNSC | NNDC | NNSR | NNDR |
| ChatGLM3-6B | 0.561 | 0.797 | 0.644 | 0.836 | 0.572 | 0.805 | 0.215 | 0.224 | 0.189 | 0.231 | 0.168 | 0.237 |
| Baichuan2-7B-Chat | 0.412 | 0.571 | 0.507 | 0.634 | 0.423 | 0.61 | 0.454 | 0.461 | 0.493 | 0.534 | 0.42 | 0.539 |
| Qwen1.5-7B-Chat | 0.742 | 0.903 | 0.835 | 0.952 | 0.803 | 0.948 | 0.503 | 0.586 | 0.606 | 0.673 | 0.526 | 0.682 |
| Llama2-7B-Chat | 0.722 | 0.81 | 0.792 | 0.857 | 0.783 | 0.845 | 0.446 | 0.429 | 0.488 | 0.513 | 0.463 | 0.494 |
| Llama3-8B-instruct | 0.77 | 0.9 | 0.891 | 0.959 | 0.868 | 0.951 | 0.644 | 0.556 | 0.725 | 0.664 | 0.707 | 0.68 |
| Baichuan2-13B-Chat | 0.414 | 0.588 | 0.484 | 0.669 | 0.409 | 0.652 | 0.309 | 0.269 | 0.336 | 0.324 | 0.303 | 0.341 |
| Qwen1.5-14B-Chat | 0.806 | 0.941 | 0.893 | 0.989 | 0.857 | 0.986 | 0.389 | 0.445 | 0.469 | 0.528 | 0.409 | 0.539 |
| Llama2-13B-Chat | 0.876 | 0.95 | 0.956 | 0.988 | 0.962 | 0.982 | 0.879 | 0.867 | 0.926 | 0.924 | 0.921 | 0.923 |
| GPT-3.5 | 0.808 | 0.862 | 0.829 | 0.92 | 0.741 | 0.898 | - | - | - | - | - | - |
| GPT-4 | 0.874 | 0.963 | 0.977 | 0.988 | 0.96 | 0.994 | - | - | - | - | - | - |
| average acc | 0.698 | 0.829 | 0.781 | 0.879 | 0.738 | 0.867 | 0.48 | 0.482 | 0.529 | 0.549 | 0.49 | 0.55 |

| Model | People Dis | | | | | | People Gen | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NSC | NDC | NNSC | NNDC | NNSR | NNDR | NSC | NDC | NNSC | NNDC | NNSR | NNDR |
| ChatGLM3-6B | 0.442 | 0.625 | 0.552 | 0.763 | 0.623 | 0.752 | 0.227 | 0.308 | 0.26 | 0.39 | 0.28 | 0.392 |
| Baichuan2-7B-Chat | 0.438 | 0.484 | 0.537 | 0.587 | 0.564 | 0.603 | 0.414 | 0.499 | 0.516 | 0.597 | 0.555 | 0.6 |
| Qwen1.5-7B-Chat | 0.634 | 0.802 | 0.805 | 0.903 | 0.876 | 0.902 | 0.504 | 0.571 | 0.632 | 0.661 | 0.701 | 0.656 |
| Llama2-7B-Chat | 0.681 | 0.706 | 0.806 | 0.834 | 0.864 | 0.819 | 0.431 | 0.494 | 0.53 | 0.594 | 0.578 | 0.598 |
| Llama3-8B-instruct | 0.675 | 0.863 | 0.831 | 0.966 | 0.888 | 0.968 | 0.572 | 0.712 | 0.695 | 0.849 | 0.739 | 0.858 |
| Baichuan2-13B-Chat | 0.473 | 0.577 | 0.551 | 0.664 | 0.569 | 0.667 | 0.316 | 0.382 | 0.385 | 0.436 | 0.418 | 0.443 |
| Qwen1.5-14B-Chat | 0.703 | 0.894 | 0.863 | 0.973 | 0.92 | 0.978 | 0.437 | 0.585 | 0.542 | 0.712 | 0.583 | 0.719 |
| Llama2-13B-Chat | 0.824 | 0.928 | 0.929 | 0.988 | 0.97 | 0.989 | 0.909 | 0.963 | 0.961 | 0.989 | 0.973 | 0.982 |
| GPT-3.5 | 0.651 | 0.707 | 0.815 | 0.851 | 0.862 | 0.849 | - | - | - | - | - | - |
| GPT-4 | 0.783 | 0.924 | 0.941 | 0.988 | 0.965 | 0.985 | - | - | - | - | - | - |
| average acc | 0.63 | 0.751 | 0.763 | 0.852 | 0.81 | 0.851 | 0.476 | 0.564 | 0.565 | 0.653 | 0.603 | 0.656 |

| Model | Place Dis | | | | | | Place Gen | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NSC | NDC | NNSC | NNDC | NNSR | NNDR | NSC | NDC | NNSC | NNDC | NNSR | NNDR |
| ChatGLM3-6B | 0.582 | 0.793 | 0.751 | 0.938 | 0.891 | 0.938 | 0.292 | 0.346 | 0.339 | 0.34 | 0.366 | 0.319 |
| Baichuan2-7B-Chat | 0.569 | 0.755 | 0.694 | 0.852 | 0.822 | 0.853 | 0.572 | 0.642 | 0.643 | 0.673 | 0.66 | 0.68 |
| Qwen1.5-7B-Chat | 0.808 | 0.906 | 0.913 | 0.979 | 0.973 | 0.976 | 0.649 | 0.714 | 0.708 | 0.75 | 0.804 | 0.76 |
| Llama2-7B-Chat | 0.745 | 0.862 | 0.839 | 0.932 | 0.928 | 0.926 | 0.423 | 0.521 | 0.508 | 0.603 | 0.56 | 0.595 |
| Llama3-8B-instruct | 0.848 | 0.93 | 0.923 | 0.979 | 0.935 | 0.979 | 0.57 | 0.666 | 0.651 | 0.793 | 0.709 | 0.808 |
| Baichuan2-13B-Chat | 0.446 | 0.637 | 0.523 | 0.823 | 0.556 | 0.819 | 0.347 | 0.431 | 0.376 | 0.498 | 0.412 | 0.504 |
| Qwen1.5-14B-Chat | 0.891 | 0.947 | 0.965 | 0.991 | 0.978 | 0.988 | 0.646 | 0.737 | 0.72 | 0.831 | 0.802 | 0.821 |
| Llama2-13B-Chat | 0.931 | 0.966 | 0.986 | 0.995 | 0.987 | 0.992 | 0.859 | 0.886 | 0.887 | 0.939 | 0.929 | 0.943 |
| GPT-3.5 | 0.799 | 0.885 | 0.872 | 0.968 | 0.909 | 0.964 | - | - | - | - | - | - |
| GPT-4 | 0.891 | 0.945 | 0.957 | 0.993 | 0.965 | 0.988 | - | - | - | - | - | - |
| average acc | 0.751 | 0.862 | 0.842 | 0.945 | 0.895 | 0.942 | 0.545 | 0.618 | 0.604 | 0.678 | 0.655 | 0.679 |

Table 10: The evaluation results for FPQs on the discriminative task (referred to as Dis) and the generative task (referred to as Gen).

| Domain | 1-hop | 2-hop | 3-hop | 4-hop | 5-hop | Total |
|--------|-------|-------|-------|-------|-------|-------|
| Art | 1988 | 342 | 500 | 748 | 1391 | 4969 |
|     | 764  | 168 | 389 | 1038 | 2610 |      |
| People | 858 | 370 | 937 | 807 | 1925 | 4897 |
|        | 923 | 234 | 575 | 943 | 2222 |      |
| Place | 237 | 244 | 610 | 1133 | 2770 | 4994 |
|       | 403 | 150 | 509 | 1043 | 2889 |      |

Table 11: For NSC and NDC, we set the distance between the edited object and the subject to one to five hops. The upper part of columns 2 to 6 presents the distribution of NSC, and the lower part shows the distribution of NDC.



Figure 11: Accuracy of all models in NSC by hops. Top: Art domain. Middle: People domain. Bottom: Place domain. Left: The discriminative task. Right: The generative task.
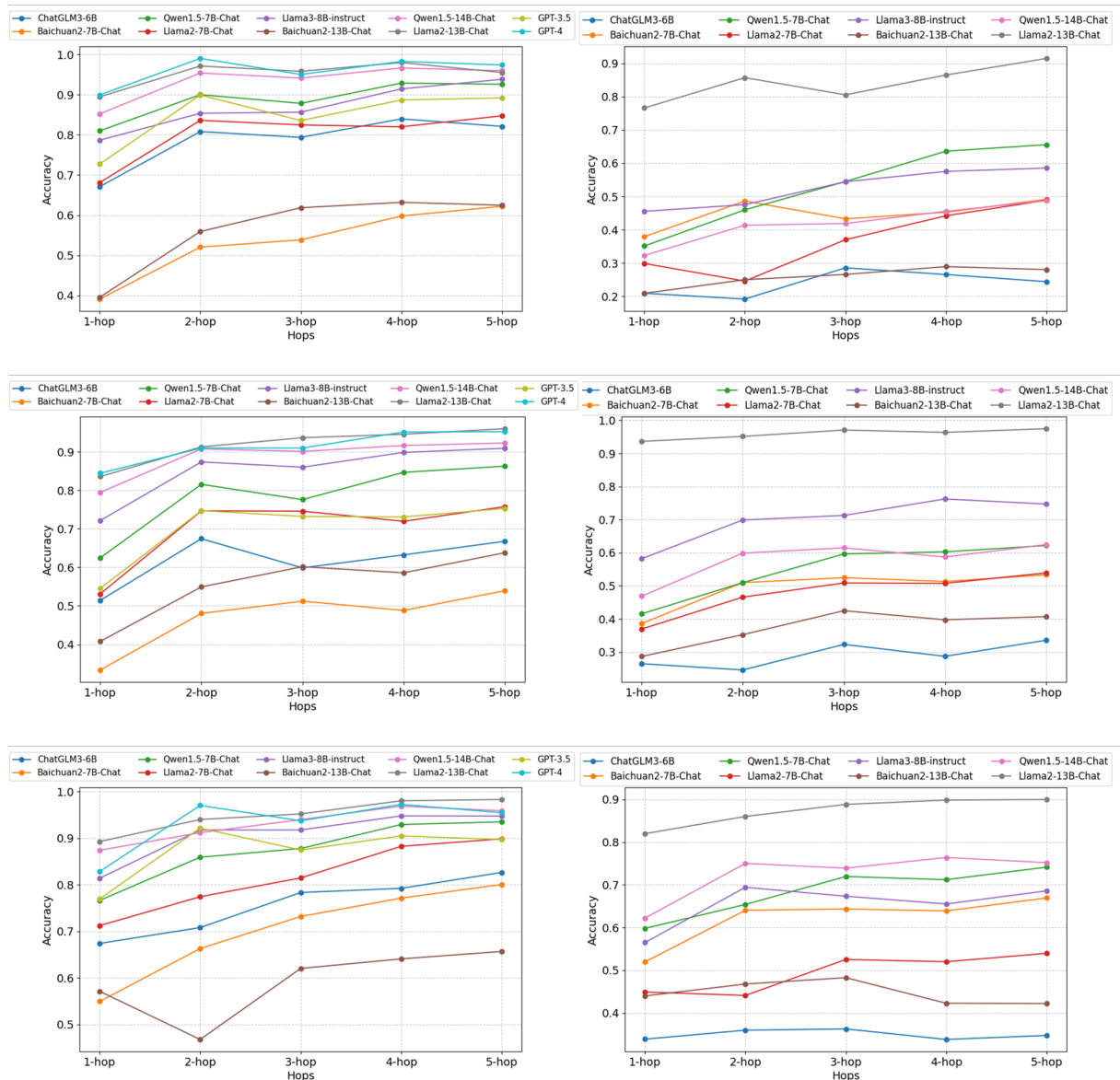
Figure 12: Accuracy of all models in NDC by hops. Top: Art domain. Middle: People domain. Bottom: Place domain. Left: The discriminative task. Right: The generative task.
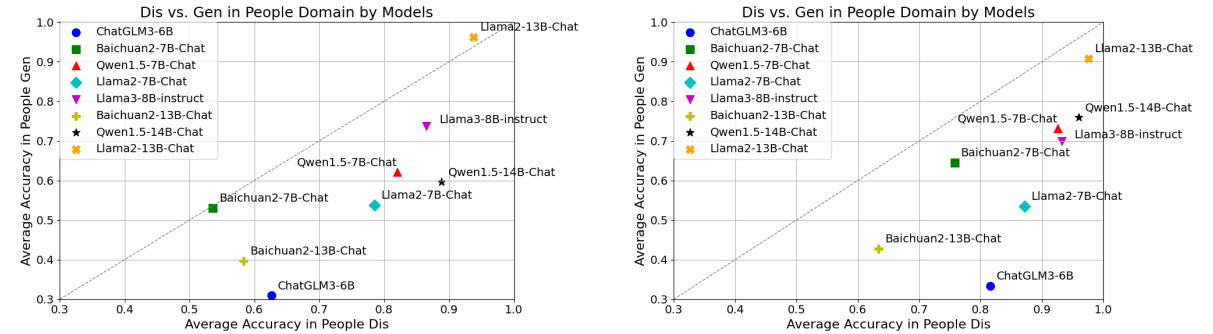


Figure 13: The overall performance comparison between the discriminative task and the generative task by models. Left: Results in People domain. Right: Results in Place domain.

| Model | Art | People | Place |
|-------|-----|--------|-------|
| ChatGLM3-6B | 0.646 | 0.752 | 0.566 |
| Baichuan2-7B-Chat | 0.892 | 0.879 | 0.596 |
| Qwen1.5-7B-Chat | 0.583 | 0.699 | 0.517 |
| Llama2-7B-Chat | 0.404 | 0.618 | 0.593 |
| Llama3-8B-instruct | 0.565 | 0.736 | 0.582 |
| Baichuan2-13B-Chat | 0.902 | 0.87 | 0.908 |
| Qwen1.5-14B-Chat | 0.563 | 0.649 | 0.444 |
| Llama2-13B-Chat | 0.191 | 0.395 | 0.343 |
| GPT-3.5 | 0.741 | 0.769 | 0.674 |
| GPT-4 | 0.649 | 0.718 | 0.632 |
| Average | 0.614 | 0.708 | 0.586 |

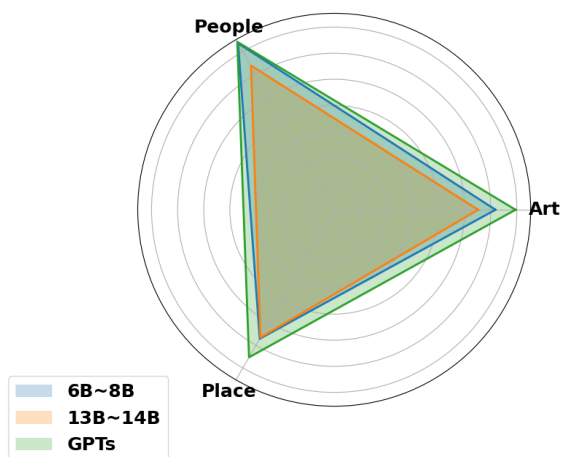Table 12: The evaluation results on Yes-No format TPQs.



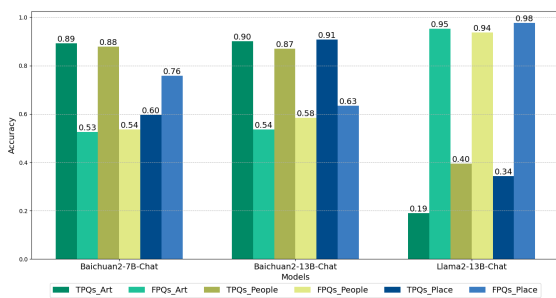Figure 14: The average accuracy of TPQs comparison across different model size.



Figure 15: The average accuracy of TPQs and FPQs across domains for Baichuan2 series and Llama2-13B-Chat.