

# Explicit Inductive Inference using Large Language Models

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) are reported to hold undesirable attestation bias on inference tasks: when asked to predict if a premise  $P$  entails a hypothesis  $H$ , instead of considering  $H$ 's conditional truthfulness entailed by  $P$ , LLMs tend to use the out-of-context truth label of  $H$  as a fragile proxy. In this paper, we propose a pipeline that exploits this bias to do explicit inductive inference. Our pipeline uses an LLM to transform a premise into a set of attested alternatives, and then aggregate answers of the derived new entailment inquiries to support the original inference prediction. On a directional predicate entailment benchmark, we demonstrate that by applying this simple pipeline, we can improve the overall performance of LLMs on inference and substantially alleviate the impact of their attestation bias.<sup>1</sup>

## 1 Introduction

Large Language Models (LLMs) are claimed to possess *implicit* inductive reasoning ability through pre-training: from the massive examples they memorized, they draw inference rules and encode them latently so that they can apply these rules to do reasoning at test time.

However, recently McKenna et al. (2023a) has pointed out that LLMs are severely affected by an attestation bias when performing inference tasks. Given the question of whether premise  $P$  entails hypothesis  $H$  with few-shot examples, an LLM's prediction is deeply bound to the hypothesis' out-of-context truthfulness, instead of its conditional truthfulness entailed by the premise. When the hypothesis  $H$  is attested in an LLM's world knowledge (the LLM believes  $H$  to be true), the LLM is likely to predict the entailment to be true, regardless of the premise. As a result, LLMs suffer a significant performance drop when the entailment labels disagree with the attestation label of hypothesis  $H$ .

<sup>1</sup>Our codes and data will be released upon publication.

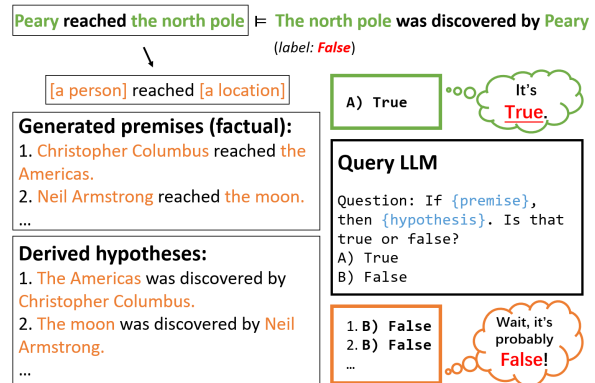


Figure 1: An example of the explicit inductive inference pipeline. While direct entailment inquiry gets a wrong answer, it can be corrected by reasoning on more alternative examples.

Although this is a severe problem limiting LLMs' performance on non-attested inferences, we argue that with careful design, this bias can instead be **exploited** to improve LLM performance on inference tasks. We notice a statistically true conclusion: Given an entailment inquiry  $P \models H$ , the attestation bias is harmful only when the premise  $P$  is not attested. If we control  $P$  to always be attested, then  $P \models H$  will naturally share the same truth label with  $H$  on a distributional basis, which dissolves the negative effects of LLMs' attestation bias.

Applying this idea, we propose a simple yet effective Explicit Inductive Inference pipeline with LLMs. As illustrated in Figure 1, the core idea is to transform a premise into a set of attested alternatives by replacing the arguments, and to aggregate the LLM's predictions on these derived inquiries to support answering the original question.

We test our pipeline with two LLMs on Levy/Holt (Levy and Dagan, 2016; Holt, 2019), a difficult directional predicate inference dataset, and further analyze the influence of our pipeline against the models' attestation bias. The results show that our pipeline can improve not only LLM's overall performance on predicate inference, but also their

robustness against the attestation bias.

To summarize our contribution, we propose an easy-to-use inference pipeline that **1)** improves LLMs’ performance on predicate inference, **2)** substantially alleviates negative effects of the LLMs’ attestation bias, and **3)** uses LLMs’ own generation capability without requiring external knowledge.

## 2 Related Work

LLMs accumulate a bias towards factual knowledge by encoding massive facts during pre-training (Roberts et al., 2020; Carlini et al., 2022; Yan et al., 2022). Recently, McKenna et al. (2023a) pointed out that LLMs suffer from an attestation bias on inference tasks as a result. Note that the effect of attestation bias is similar to that of the hypothesis-only baseline (Poliak et al., 2018), but while the former is a bias from pre-training, the latter originates from dataset artifacts in supervised learning.

In other tasks, previous work has mitigated the bias towards attestation by introducing counterfactual examples (Wang et al., 2022b; Zhou et al., 2023; Wang et al., 2023) or replacing argument entities with their type labels (Zhou et al., 2024). In this paper, we go one step further to show that in an inference task, we should instead encourage the models to generate factual alternative examples.

The idea of aggregating multiple versions of LLMs’ outputs has been explored in prior work. Wang et al. (2022a) encourage LLMs to generate multiple reasoning paths for one question, while Zhou et al. (2022) embody one question with multiple prompts. In contrast, our method creates semantically different alternative questions, which serve as extra evidence for an original inquiry.

## 3 Explicit Inductive Inference

### 3.1 Task and Definition

The task of this work is to determine the entailment relation between two binary predicates where both predicates are contextualized with the same pair of entities. The input will be in the form of two triples  $(s, p, o) - (s, h, o)$  where  $s$  is the subject entity,  $o$  is the object entity,  $p$  is the premise predicate, and  $h$  is the hypothesis predicate. There are also cases in the form of  $(s, p, o) - (o, h, s)$  where the two entities are swapped in position like the example in Figure 1. Without loss of generality, we describe our method with inputs in the former format.

The goal is to predict whether the premise triple entails the hypothesis triple, namely the truth label

of  $(s, p, o) \models (s, h, o)$ . To use an LLM to predict entailments, each input triple pair will be wrapped into a prompt. We mark them as  $Q[(s, p, o) \models (s, h, o)]$  and call them entailment inquiries.

### 3.2 Exploit the Attestation Bias

As stated in Section 1, the attestation bias of LLMs can be less detrimental if the premise  $P$  is attested in an entailment inquiry, because the truth label of  $P \models H$  would likely be the same as the attestation label of  $H$ . Besides this, two more insights are guiding our pipeline design:

1) The label of a predicate entailment inquiry does not change when the argument entities are replaced, as long as the substitution entities keep the same semantic type labels.

2) Factual  $\neq$  Attested. Factual knowledge from external sources may not be confirmed by LLMs for being longtail, absent in pre-training data, or conflicted with out-of-date records. Facts generated by LLMs, on the other hand, are highly likely to be recognizable by themselves. Even hallucinated generations are acceptable since they are still attested if not factual.

Based on these understandings, we propose the **Explicit InDuctive Inference (EIDI)** pipeline. Given an entailment inquiry  $P \models H$ , our EIDI pipeline first transforms  $P$  into a set of different attested premises  $P'$ s by replacing the arguments of  $P$ . Then the corresponding set of  $H'$ s is derived, so that we now have a list of alternative inquiries  $P' \models H'$ . Finally, we explicitly do an inductive inference on these new inquiries by drawing a concluded entailment prediction from an LLM’s answers to these alternative inquiries.

It is worth mentioning that given  $P$  is true, logically,  $H$  is always true if  $P \models H$  but not vice versa. We can only statistically conjecture  $P \models H$  if we observe a high probability of  $H$  being true (predicted by the LLM according to the bias). Therefore, we encourage the transformation module to generate a variety of different alternative premise triples, so that a more reliable conclusion can be drawn when we aggregate the predictions.

### 3.3 Explicit Inductive Inference Pipeline

**Typing** While the label of (medicine X, kills, disease Y)  $\models$  (medicine X, is a cure of, disease Y) is *True*, one can not therefore deduce that (Person X, kills, animal Y)  $\models$  (Person X, is a cure of, Animal Y). To prevent these errors incited by the ambiguity of predicates, for each premise triple  $(s, p, o)$ , we

164 query the LLM to obtain the entity type label of the  
165 subject and object  $t_s$  and  $t_o$ . Here we do not prede-  
166 fine a vocabulary for possible type labels since the  
167 purpose is only to disambiguate.

168 **Transformation** With these assigned type labels we query the LLM to generate alternative  
169 arguments for the premise predicate. From one  
170 typed premise triple  $(s, t_s, p, o, t_o)$ , we encour-  
171 age the LLM to generate a list of new attested  
172 triples  $(s_1, p, o_1), \dots, (s_n, p, o_n)$  where the substi-  
173 tution entities keep the original types, i.e. any  
174  $s_i$  still has type  $t_s$  and any  $o_i$  still has type  $t_o$ .  
175 These  $n$  new premise triples will then be expanded  
176 to  $n$  new entailment inquiries  $Q[(s_1, p, o_1) \models$   
177  $(s_1, h, o_1)], \dots, Q[(s_n, p, o_n) \models (s_n, h, o_n)]$ .

179 **Prediction** At this point, we input each derived  
180 entailment inquiry  $Q[(s_i, p, o_i) \models (s_i, h, o_i)]$  to  
181 the LLM to get their response and corresponding  
182 probability score. Then we take the average score  
183 of these  $n$  different scores as our explicit inductive  
184 score for the original entailment inquiry.

## 185 4 Experimental Setup

### 186 4.1 Datasets

187 We test our pipeline on the Levy/Holt dataset  
188 (Levy and Dagan, 2016; Holt, 2019), a predicate  
189 entailment dataset where each entry consists of  
190 two triples in the form of  $(s, p, o) - (s, h, o)$  or  
191  $(s, p, o) - (o, h, s)$ , and a following label shows  
192 whether the premise triple entails the hypothesis  
193 triple. We use the directional portion of this dataset  
194 following prior work (McKenna et al., 2023b; Chen  
195 et al., 2022; Li et al., 2022), as it is a challenging  
196 test focused on the understanding of entailment  
197 beyond bi-directional similarity.

198 Following McKenna et al. (2023a), we further  
199 analyze how the LLMs’ attestation bias is digested  
200 in our method. We split the Levy/Holt dataset ac-  
201 cording to whether the label of  $P \models H$  agrees  
202 with the attestation label (obtained by querying  
203 the LLM) of  $H$  for each entry. For the 1784 en-  
204 tries in the full directional test set, this yields an  
205 attestation-consistent subset of 956 entries and an  
206 attestation-adversarial subset of 828 entries.<sup>2</sup> We  
207 report results on both the directional test set and its  
208 two subsets in Section 5.

<sup>2</sup>The substantial size of the attestation-adversarial subset demonstrates the detrimental effect of attestation bias in real datasets.

### 209 4.2 LLMs

210 We test our method with two LLMs, GPT-3.5 and  
211 Llama3. GPT-3.5 (OpenAI, 2023) is a set of power-  
212 ful closed-source commercial LLMs. We choose  
213 the GPT-3.5-Turbo version for its widespread use  
214 in the research community. Llama3 (Meta, 2024)  
215 is a SOTA open-source LLM, where we choose  
216 the largest Llama3-70B-instruct version for its op-  
217 timized capacity. Throughout our experiments, we  
218 use greedy decoding for reproducible results.

219 Our pilot studies on the development set indi-  
220 cate that adding few-shot examples in the predic-  
221 tion module may add extra bias to the model, and  
222 therefore introduce unnecessary considerations on  
223 finding proper examples under each setting. Hence  
224 we choose zero-shot prompts for the prediction  
225 module and one-shot prompts for the transforma-  
226 tion module where the only example is the original  
227 premise. Discussion on prompt selection and a list  
228 of all prompts we use are included in Appendix A.

### 229 4.3 Baselines and Metric

230 We compare EIDI against two baselines. We con-  
231 struct  $MCQ_{entity}$  baseline by directly wrapping the  
232 original premise and hypothesis with the Multiple-  
233 Choice Question prompt used in our prediction  
234 module, and passing it to the LLM to get an en-  
235 tailment prediction.  $MCQ_{type}$  baseline is set up  
236 in the same way where the only difference is that  
237 we first replace the arguments of the predicates by  
238 their entity types. To keep ourselves aligned with  
239 previous work, we use the 48 FIGER types (Ling  
240 and Weld, 2012) as in McKenna et al. (2023a) for  
241 this measure, instead of the LLM-generated types  
242 in Section 3.3.

243 We draw the precision-recall curve for EIDI and  
244 each baseline by inspecting the final output token  
245 probability of the model’s response. As a result  
246 of the multiple-choice prompt design, returned an-  
247 swers always start with a choice mark where A is  
248 for entailment and B is for non-entailment. For  
249 baseline methods, we score that one token’s prob-  
250 ability. For our EIDI pipeline, we calculate the  
251 average score of the  $k$  output tokens’ probabilities.

252 Following Li et al. (2022); McKenna et al.  
253 (2023a), we calculate the normalized area-under-  
254 curve ( $AUC_{norm}$ ) as an indicator of the model’s  
255 performance. This measure describes how much  
256 better a model is over a degenerate baseline return-  
257 ing positive answers to every data entry.

Pipeline	Model	
	GPT-3.5	Llama3
MCQ <sub>entity</sub>	23.85	36.66
MCQ <sub>type</sub>	25.88	35.13
EIDI <sub>all</sub>	<b>35.52</b>	<b>50.89</b>
EIDI <sub>1</sub>	31.16	41.85
EIDI <sub>2</sub>	32.10	46.75
EIDI <sub>5</sub>	33.41	49.61

Table 1: Overall normalized Area-Under-the-Curve (%) of our EIDI pipeline and the two baselines on the full Levy/Holt directional test set. EIDI<sub>*i*</sub> inspects only *i* alternative inquiries, and EIDI<sub>all</sub> considers all examples obtained in the transformation step.

## 5 Results and Discussion

### 5.1 Overall performance

Table 1 shows the performance of each model on the directional Levy/Holt test set. With both LLMs, our EIDI<sub>all</sub> pipeline gains a significant improvement over the two baseline methods.

The typical value of the size of total generated examples  $n$  is 10 for the EIDI<sub>all</sub> setting. It can be observed that the performance of EIDI<sub>*i*</sub> steadily increases along with *i*, confirming our hypothesis that with attested  $P'$ s, the more cases of alternative  $P' \models H'$  generated, the more reliable our pipeline is. The complete results of all EIDI<sub>*i*</sub> settings are shown in Appendix B.

An interesting observation lies between the performance of the EIDI<sub>1</sub> setting and the baselines, which shows that replacing the original inquiry with even one self-generated example can improve the LLMs’ predicate inference performance. The difference between EIDI<sub>1</sub> and MCQ<sub>type</sub> baseline also highlights the importance of instantiating attested triples. Since the effect of the attestation bias is already excluded from the results of the MCQ<sub>type</sub>, this proves that the EIDI pipeline is taking advantage of further exploiting the bias.

### 5.2 Against the bias

Table 2 compares the performance of each method on attestation-consistent (*cons.*) and attestation-adversarial (*adv.*) subsets. Measured by the difference of AUC<sub>norm</sub> between the two subsets, our pipeline reduces the effect of LLMs’ attestation bias by over 20% from the MCQ<sub>type</sub> baseline, and over 35% from the MCQ<sub>entity</sub> baseline in average.

With both LLMs, we observe an AUC<sub>norm</sub> of near 0% in the two baseline settings, demonstrating

Model	Pipeline	<i>cons.</i>	<i>adv.</i>	<i>diff.</i>
GPT-3.5	MCQ <sub>entity</sub>	82.04	0.00	-82.04
	MCQ <sub>type</sub>	69.40	0.48	-68.92
	EIDI <sub>all</sub>	56.14	9.97	<b>-46.17</b>
	EIDI <sub>1</sub>	53.73	8.95	<b>-44.78</b>
Llama3	MCQ <sub>entity</sub>	81.08	0.01	-81.07
	MCQ <sub>type</sub>	70.25	2.41	-67.84
	EIDI <sub>all</sub>	69.59	23.83	<b>-45.76</b>
	EIDI <sub>1</sub>	63.98	15.66	-48.32

Table 2: AUC<sub>norm</sub> (%) on the attestation-bias-split datasets. The *diff.* column marks the difference between results on the attestation-consistent (*cons.*) and attestation-adversarial (*adv.*) subsets.

the extreme inability of the LLMs to capture the essential entailment signal against the attestation bias in a zero-shot setting.

Interesting results appear again under the EIDI<sub>1</sub> setting. On GPT-3.5-turbo, it slightly outperforms the EIDI<sub>all</sub> setting. But this only happens because EIDI<sub>all</sub> setting is doing better on the attestation-consistent subset, which implies that EIDI<sub>all</sub> setting is still the choice for best performance, while EIDI<sub>1</sub> is also a strong candidate for scenarios with limited compute.

These results suggest that our pipeline can be used to improve LLMs’ general inference performance, and especially in attestation-adversarial scenarios, e.g. *If lions are fed on hay, then lions eat hay*. As a replacement to LLM’s direct inference prediction, EIDI pipeline can be easily plugged into frameworks with LLMs to do various downstream tasks like question answering and KG completion.

## 6 Conclusions

We propose an explicit inductive pipeline exploiting the attestation bias of LLMs to do more robust predicate inference. With experiments on the directional Levy/Holt dataset and its attestation-bias-split subsets, we have shown that our baseline gains a significant improvement over LLM’s primary inference performance, and substantially reduces the performance loss caused by LLMs’ attestation bias.

Without external knowledge, EIDI use LLMs’ own generation to exploit their attestation bias. Our results suggest that although biases of LLMs are usually undesirable obstacles, in some scenarios they may be tapped for good with careful design. We advocate for similar ideas to be applied to other tasks to improve LLM performance in future work.

## 328 Limitations

329 In this paper, we demonstrated the performance  
330 of our pipeline by comparing it to two baselines.  
331 Although we intend to exclude prompt engineering  
332 factors from our analysis, it has been widely ac-  
333 cepted that including few-shot examples and other  
334 prompting techniques can guide LLMs to output  
335 better answers. Therefore there could be further  
336 studies on evaluating the effects of using different  
337 prompts in the EIDI pipeline.

338 Generating alternative inquiries and respectively  
339 doing inferences on them can be computationally  
340 expensive compared to only one determination in  
341 baseline settings. As a result, downstream applica-  
342 tions may find a trade-off between computational  
343 efficiency and better inference performance.

344 We also tested our pipeline against the frequency  
345 bias that McKenna et al. (2023a) pointed out, and  
346 the results show that the EIDI pipeline amplifies  
347 this bias compared to the baselines due to its choice  
348 of popular entities. We argue that this reaffirms the  
349 challenge in achieving Pareto improvements on  
350 LLM robustness against biases, and leave those  
351 results and discussions to Appendix C.

## 352 References

353 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,  
354 Katherine Lee, Florian Tramer, and Chiyuan Zhang.  
355 2022. Quantifying memorization across neural lan-  
356 guage models. *arXiv preprint arXiv:2202.07646*.

357 Zhibin Chen, Yansong Feng, and Dongyan Zhao. 2022.  
358 Entailment graph learning with textual entailment  
359 and soft transitivity. In *Proceedings of the 60th An-  
360 nual Meeting of the Association for Computational  
361 Linguistics (Volume 1: Long Papers)*, pages 5899–  
362 5910, Dublin, Ireland. Association for Computational  
363 Linguistics.

364 Xavier Holt. 2019. Probabilistic models of relational  
365 implication. *Preprint*, arXiv:1907.12048.

366 Omer Levy and Ido Dagan. 2016. Annotating rela-  
367 tion inference in context via question answering. In  
368 *Proceedings of the 54th Annual Meeting of the As-  
369 sociation for Computational Linguistics (Volume 2:  
370 Short Papers)*, pages 249–255, Berlin, Germany. As-  
371 sociation for Computational Linguistics.

372 Tianyi Li, Mohammad Javad Hosseini, Sabine Weber,  
373 and Mark Steedman. 2022. Language models are  
374 poor learners of directional inference. In *Findings  
375 of the Association for Computational Linguistics:  
376 EMNLP 2022*, pages 903–921, Abu Dhabi, United  
377 Arab Emirates. Association for Computational Lin-  
378 guistics.

Xiao Ling and Daniel Weld. 2012. Fine-grained entity  
recognition. In *Proceedings of the AAAI Conference  
on Artificial Intelligence*, volume 26, pages 94–100. 379  
380  
381

Nick McKenna, Tianyi Li, Liang Cheng, Mohammad  
Hosseini, Mark Johnson, and Mark Steedman. 2023a.  
[Sources of hallucination by large language models  
on inference tasks](#). In *Findings of the Association  
for Computational Linguistics: EMNLP 2023*, pages  
2758–2774, Singapore. Association for Computa-  
tional Linguistics. 382  
383  
384  
385  
386  
387  
388

Nick McKenna, Tianyi Li, Mark Johnson, and Mark  
Steedman. 2023b. [Smoothing entailment graphs with  
language models](#). In *Proceedings of the 13th Inter-  
national Joint Conference on Natural Language Pro-  
cessing and the 3rd Conference of the Asia-Pacific  
Chapter of the Association for Computational Lin-  
guistics (Volume 1: Long Papers)*, pages 551–563,  
Nusa Dua, Bali. Association for Computational Lin-  
guistics. 389  
390  
391  
392  
393  
394  
395  
396  
397

Meta. 2024. [Llama3](#). 398

OpenAI. 2023. [Openai](#). 399

Adam Poliak, Jason Naradowsky, Aparajita Haldar,  
Rachel Rudinger, and Benjamin Van Durme. 2018.  
[Hypothesis only baselines in natural language infer-  
ence](#). In *Proceedings of the Seventh Joint Confer-  
ence on Lexical and Computational Semantics*, pages  
180–191, New Orleans, Louisiana. Association for  
Computational Linguistics. 400  
401  
402  
403  
404  
405  
406

Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.  
[How much knowledge can you pack into the param-  
eters of a language model?](#) In *Proceedings of the  
2020 Conference on Empirical Methods in Natural  
Language Processing (EMNLP)*, pages 5418–5426,  
Online. Association for Computational Linguistics. 407  
408  
409  
410  
411  
412

Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou,  
and Muhao Chen. 2023. [A causal view of entity  
bias in \(large\) language models](#). In *Findings of the  
Association for Computational Linguistics: EMNLP  
2023*, pages 15173–15184, Singapore. Association  
for Computational Linguistics. 413  
414  
415  
416  
417  
418

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc  
Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery,  
and Denny Zhou. 2022a. Self-consistency improves  
chain of thought reasoning in language models. *arXiv  
preprint arXiv:2203.11171*. 419  
420  
421  
422  
423

Yiwei Wang, Muhao Chen, Wenxuan Zhou, Yujun  
Cai, Yuxuan Liang, Dayiheng Liu, Baosong Yang,  
Juncheng Liu, and Bryan Hooi. 2022b. [Should we  
rely on entity mentions for relation extraction? debi-  
asing relation extraction with counterfactual analysis](#).  
In *Proceedings of the 2022 Conference of the North  
American Chapter of the Association for Computa-  
tional Linguistics: Human Language Technologies*,  
pages 3071–3081, Seattle, United States. Association  
for Computational Linguistics. 424  
425  
426  
427  
428  
429  
430  
431  
432  
433

- 434 Jun Yan, Yang Xiao, Sagnik Mukherjee, Bill Yuchen  
435 Lin, Robin Jia, and Xiang Ren. 2022. [On the ro-](#)  
436 [bustness of reading comprehension models to entity](#)  
437 [renaming](#). In *Proceedings of the 2022 Conference*  
438 *of the North American Chapter of the Association*  
439 *for Computational Linguistics: Human Language*  
440 *Technologies*, pages 508–520, Seattle, United States.  
441 Association for Computational Linguistics.
- 442 Ben Zhou, Hongming Zhang, Sihao Chen, Dian Yu,  
443 Hongwei Wang, Baolin Peng, Dan Roth, and Dong  
444 Yu. 2024. Conceptual and unbiased reasoning in  
445 language models. *arXiv preprint arXiv:2404.00205*.
- 446 Chunting Zhou, Junxian He, Xuezhe Ma, Taylor Berg-  
447 Kirkpatrick, and Graham Neubig. 2022. [Prompt con-](#)  
448 [sistency for zero-shot task generalization](#). In *Find-*  
449 *ings of the Association for Computational Linguistics:*  
450 *EMNLP 2022*, pages 2613–2626, Abu Dhabi, United  
451 Arab Emirates. Association for Computational Lin-  
452 guistics.
- 453 Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and  
454 Muhao Chen. 2023. [Context-faithful prompting](#)  
455 [for large language models](#). In *Findings of the As-*  
456 *sociation for Computational Linguistics: EMNLP*  
457 *2023*, pages 14544–14556, Singapore. Association  
458 for Computational Linguistics.

## A Prompts Selection

Here we list all the prompts that we use in our experiments.

**Typing** The purpose of this module is only to disambiguate the predicates, therefore no vocabulary of allowed type labels is predefined.

Type the entities in the following triples:

Hitler | was born in | Poland -> a person | was born in | a country

Hogs | eats | Corn -> an animal | eats | a food

Aspirin | may reduce the risk of | Cancer -> a medicine | may reduce the risk of | a disease

$\{s\} | \{p\} | \{o\} \rightarrow$

**Transformation** Although we use the word 'fact', the generated triples are always attested rather than factual.

Write  $\{n + 1\}$  facts in the form of " $\{t_s\} | \{p\} | \{t_o\}$ ."

-  $\{s\} | \{p\} | \{o\}$ .

-

**Prediction** This is also used for the two baselines.

Question: If  $\{s\} | \{p\} | \{o\}$ , then  $\{s\} | \{h\} | \{o\}$ . Is that true or false?

Choices:

A) True

B) False

Answer:

For prediction module, when an instruction is required, we use the following one:

Only return one mark A, B or C to answer the question.

## B Results on all EIDI<sub>i</sub> Settings

Table 3 shows the performance of all EIDI<sub>i</sub> settings. Best performances are reached when all transformed alternative inquiries are considered.

Pipeline	Model	
	GPT-3.5	Llama3
MCQ <sub>entity</sub>	23.85	36.66
MCQ <sub>type</sub>	25.88	35.13
EIDI <sub>1</sub>	31.16	41.85
EIDI <sub>2</sub>	32.10	46.75
EIDI <sub>3</sub>	31.47	47.52
EIDI <sub>4</sub>	32.05	48.60
EIDI <sub>5</sub>	33.54	49.61
EIDI <sub>6</sub>	33.41	50.42
EIDI <sub>7</sub>	34.68	50.13
EIDI <sub>8</sub>	34.76	50.36
EIDI <sub>9</sub>	35.28	50.39
EIDI <sub>10</sub>	<b>35.52</b>	50.01
EIDI <sub>11</sub>	-	50.52
EIDI <sub>12</sub>	-	<b>50.89</b>

Table 3: AUC<sub>norm</sub> (%) of all EIDI<sub>i</sub> settings.

## C Frequency Bias

We also tested our pipeline on the frequency bias using the same dataset split measure as that for attestation bias. The dataset that we use is from McKenna et al. (2023a)'s work, where we have 972 entries of frequency-consistent input and 220 entries of frequency-adversarial input.

Compared to baselines, the EIDI pipeline introduces extra frequency bias. This is expected since our transformation module is not designed to alter the relative frequency of the predicates, and may have amplified the frequency bias by taking popular alternative entities generated by the LLMs. This result reaffirms the challenging nature of directional inference and the difficulty in improving robustness against multiple biases at once.

Model	Pipeline	cons.	adv.	diff.
GPT-3.5	MCQ <sub>entity</sub>	20.58	29.38	+8.80
	MCQ <sub>type</sub>	24.49	32.93	+8.44
	EIDI <sub>all</sub>	40.66	20.83	-19.83
	EIDI <sub>1</sub>	33.94	18.83	-15.11
Llama3	MCQ <sub>entity</sub>	33.30	47.87	+14.57
	MCQ <sub>type</sub>	31.47	47.19	+15.72
	EIDI <sub>all</sub>	51.97	42.27	-9.70
	EIDI <sub>1</sub>	39.78	35.32	-4.46

Table 4: Normalized area-under-curve(%) on the frequency-bias-split datasets.

## **D Computational Cost**

Our experiments on Llama3-70B-Instruct are applied on two A6000 GPUs. For 1784 entries and 10 alternative inquiries for each entry, the typing module takes about 3 GPU hour, the transformation module takes about 100 GPU hours, and the prediction module takes about 6 GPU hours.