

# When LLMs Meet Cunning Texts: A Fallacy Understanding Benchmark for Large Language Models

Yinghui Li<sup>1</sup>; Qingyu Zhou<sup>2,\*</sup>; Yuanzhen Luo\*, Shirong Ma<sup>1</sup>,  
 Yangning Li<sup>1</sup>, Hai-Tao Zheng<sup>1,†</sup>, Xuming Hu<sup>3,†</sup>, Philip S. Yu<sup>4</sup>  
<sup>1</sup>Tsinghua University, <sup>2</sup>ByteDance Inc.  
<sup>3</sup>The Hong Kong University of Science and Technology (Guangzhou)  
<sup>4</sup>University of Illinois Chicago  
 liyinghu20@mails.tsinghua.edu.cn



## Abstract

Recently, Large Language Models (LLMs) make remarkable evolutions in language understanding and generation. Following this, various benchmarks for measuring all kinds of capabilities of LLMs have sprung up. In this paper, we challenge the reasoning and understanding abilities of LLMs by proposing a **FaLLacy Understanding Benchmark (FLUB)** containing cunning texts that are easy for humans to understand but difficult for models to grasp. Specifically, the cunning texts that FLUB focuses on mainly consist of the tricky, humorous, and misleading texts collected from the real internet environment. And we design three tasks with increasing difficulty in the FLUB benchmark to evaluate the fallacy understanding ability of LLMs. Based on FLUB, we investigate the performance of multiple representative and advanced LLMs, reflecting our FLUB is challenging and worthy of more future study. Interesting discoveries and valuable insights are achieved in our extensive experiments and detailed analyses. We hope that our benchmark can encourage the community to improve LLMs' ability to understand fallacies. Our data and codes are available at <https://github.com/THUKEI/FLUB>.

## The FLUB Benchmark

we collect real cunning texts as our raw data from a famous Chinese online forum, the "Ruozhiba". This forum is popular for its cunning and unreasonable posts, which are generally easy for humans to understand but challenging for LLMs. The characteristics of the posts contained in this forum are consistent with our research motivation, so choosing it as the data source well supports FLUB's evaluation of LLMs' fallacy understanding ability. After data cleaning and annotating of cunning types, FLUB has 8 fine-grained types of cunning texts and most of the texts in FLUB fall into two types of fallacy, namely, faulty reasoning and word game. Moreover, we also manually annotated one correct answer (i.e., the explanation of the cunning text) and three confusing wrong answers for each input text in FLUB.

## The Cunning Types of FLUB

Cunning Type	Definition	Example
错误类比 False Analogy	由于事件A和事件B具有某种相似性，从而得出事件A和事件B也具有其他属性。或者错误地比出事件A和事件B具有其他属性。 Due to the occurrence of event A being or being accompanied by certain attribute, it is erroneously analogized that event B, which is similar to event A, should also have that attribute, or that event B, which is opposite to event A, should have the opposite attribute.	很多人担心出门忘记关门。为什么?因为门不会自己开门! Many people worry about forgetting to close the door when they leave home. Why don't they worry about whether they have opened the door when they come in?
冷笑话 Lame Jokes	由于缺乏对某个常识或事实的认知，从而得出某个不符合逻辑的结论。 Due to a lack of understanding of a common sense or fact, a illogical question or conclusion can be drawn. Note that this sentence may be funny due to its unusual logical implication.	忘记把钱存在哪个ATM机里了怎么办? 银行不会记得吗, 还是都记不住? What should I do if I forget which ATM I deposited money into? The bank has several ATMs, and they all look the same.
字音错误 Phonetic Error	注意, 该句子包含故意改变字音的词语以制造歧义或误导读者。 Note that this sentence may be funny due to its unusual logical implication. 通过改变汉语中某些字的发音, 从而制造出歧义或误导读者的句子。 By changing the meaning of a polysemous word in a sentence, illogical questions or conclusions can be drawn.	语文老师问我中午吃什么, 我说我带了个三明治。 My teacher said the sentence is grammatically wrong. Should I give this sentence some articulation or administer an IV drip?
歧义 Ambiguity	通过改变句子的某个词义, 从而得出与原文不符的结论。 By changing the meaning of a polysemous word in a sentence, illogical questions or conclusions can be drawn.	语文老师问我中午吃什么, 我说我带了个三明治。 My teacher said the sentence is grammatically wrong. Should I give this sentence some articulation or administer an IV drip?
悖论 Paradox	句子的意思和表述互相矛盾。 The expression of a sentence or question is contradictory.	"凡属绝对者皆是绝对者" Is too absolute.
事实性错误 Factual Error	由于缺乏对某个事实的认知, 或者对事实的认知, 从而得出与原文不符的结论。 Due to a lack of understanding or distortion of fact, meaningless questions or conclusions are raised.	一块铁和一块棉花哪个重? Which one weighs more, a ton of iron or a ton of cotton?
推理错误 Reasoning Error	从一事件中推导出一个错误或不相关的结论, 或者忽略了事件间的因果关系。 Inferring an incorrect or meaningless conclusion from an event, or reversing the causal relationship of the event.	根据我在非洲国家看到的, 我的人口老龄化已经相当严重了。 According to my sorry tale from traveling homes, the aging of the population in our country has become quite severe.
文字游戏 Word Game	错误地改变中文文字的意思或用法, 在基础上提出逻辑问题或得出结论。 Mistakenly changing the meaning of words in a sentence, and based on this, raising questions or drawing conclusions.	人类70%是水, 所以10个人有7个人是水的成分! 70% of the human body is water, so 7 out of 10 people are water disguised as humans!
未分类 Undefined	句子本身具有错误, 或者句子的表述不符合逻辑, 但是不属于上述任何一个类别。 The sentence itself has errors, or the expression of the sentence does not conform to normal logic, but does not belong to any of the above categories.	在逻辑上能成立并能使用吗? Is it feasible to open a bar at a highway service area?

## The Benchmark Task Setups of FLUB

- Task 1: Answer Selection.** In Task 1, LLMs are required to select the correct answer from four given candidate explanations for each input text. The design motivation of this task is to test whether LLMs can distinguish right from wrong when seeing the correct and wrong answers in the context of a given cunning text.
- Task 2: Cunning Type Classification.** If LLMs are directly tasked with determining the corresponding cunning type, it will help us in conducting an initial automated assessment of the LLM's understanding ability. The cunning type classification task is specifically designed to evaluate whether LLMs can classify the cunning text into categories aligned with human intuition based on the hidden irrational aspects within the current text.
- Task 3: Fallacy Explanation.** To further test whether LLMs truly understand the given cunning text, we design the explanation task. In this task, the designed prompt and input texts are directly input into LLMs, enabling them to "read" input texts and generate corresponding explanations.
- Automatic Evaluation Metrics:** For Task 1, we calculate **Accuracy** directly based on the LLMs' selection results. For Task 2, considering that there are a few cunning types in FLUB with small sample size, we choose the **F-1 Score** to measure the performance of LLMs. For Task 3, we assign a **GPT-4 Score** ranging from 1 to 10.
- Human Evaluation Settings:** For Task 1 and Task 2, we conduct human evaluations to explore how well human-level intelligence could perform these two tasks. For the human evaluation of Task 3, we mainly want to verify the effectiveness of the automatic GPT-4 score we use, therefore, we hire 3 evaluation annotators to rate LLMs' explanations, with scores ranging from {1, 2, 3, 4, 5}.

## Automatic Evaluation Results

Table 1: We bold the optimal and underline the suboptimal of closed/open-source models. We report the overall performance by calculating the **geometric mean** of the three tasks. We color the result that Chain-of-Thought (CoT) brings positive/negative gain as green<sup>↑</sup>/red<sup>↓</sup>.

Models	Open Source	Selection Accuracy		Classification F-1 Score		Explanation GPT-4 Score		Overall Performance	
		w/o CoT	CoT	w/o CoT	CoT	w/o CoT	CoT	w/o CoT	CoT
ERNIE-Bot-3.5-Turbo [15]	X	32.97	34.65 <sup>↑</sup>	1.99	6.09 <sup>↑</sup>	5.78	5.83 <sup>↑</sup>	7.24	10.72 <sup>↑</sup>
ERNIE-Bot-4.0 [15]	X	52.76	38.37 <sup>↓</sup>	10.33	11.15 <sup>↑</sup>	6.35	6.22 <sup>↓</sup>	15.13	13.86 <sup>↓</sup>
ERNIE-Bot-4.0 [15]	X	75.66	71.34 <sup>↓</sup>	11.84	14.42 <sup>↑</sup>	7.73	8.11 <sup>↑</sup>	19.06	20.28 <sup>↑</sup>
GPT-3.5-Turbo [16]	X	50.48	48.08 <sup>↓</sup>	3.09	6.15 <sup>↑</sup>	6.23	7.00 <sup>↑</sup>	9.91	12.74 <sup>↑</sup>
GPT-4-Turbo [16]	X	79.38	82.73 <sup>↑</sup>	12.31	13.97 <sup>↑</sup>	8.95	9.21 <sup>↑</sup>	20.60	22.00 <sup>↑</sup>
ChatGLM3-6B [17]	✓	35.85	35.01 <sup>↓</sup>	7.48	9.34 <sup>↑</sup>	4.98	4.82 <sup>↓</sup>	11.01	11.64 <sup>↑</sup>
Qwen-72B-Chat [18]	✓	38.49	33.69 <sup>↓</sup>	8.00	10.97 <sup>↑</sup>	5.39	5.65 <sup>↑</sup>	11.84	11.98 <sup>↑</sup>
Qwen-14B-Chat [18]	✓	42.57	43.05 <sup>↑</sup>	10.34	10.44 <sup>↑</sup>	5.24	6.24 <sup>↑</sup>	13.21	14.10 <sup>↑</sup>
Qwen-72B-Chat [18]	✓	58.63	61.51 <sup>↑</sup>	9.32	12.26 <sup>↑</sup>	7.34	7.90 <sup>↑</sup>	15.89	18.13 <sup>↑</sup>
Yi-6B-Chat [19]	✓	32.37	29.26 <sup>↓</sup>	8.87	9.84 <sup>↑</sup>	5.73	5.39 <sup>↓</sup>	11.81	11.58 <sup>↓</sup>
Yi-34B-Chat [19]	✓	47.96	48.80 <sup>↑</sup>	4.74	11.70 <sup>↑</sup>	6.97	7.52 <sup>↑</sup>	11.66	16.17 <sup>↑</sup>
Baichuan2-7B-Chat [20]	✓	43.17	37.17 <sup>↓</sup>	1.02	4.45 <sup>↓</sup>	5.48	4.85 <sup>↓</sup>	6.23	9.29 <sup>↓</sup>
Baichuan2-13B-Chat [20]	✓	37.05	38.01 <sup>↑</sup>	3.52	4.58 <sup>↓</sup>	5.79	5.84 <sup>↑</sup>	9.11	10.06 <sup>↑</sup>
Random	-	-	25.00	-	7.90	-	-	-	-
Human	-	-	93.35	-	63.69	-	-	-	-

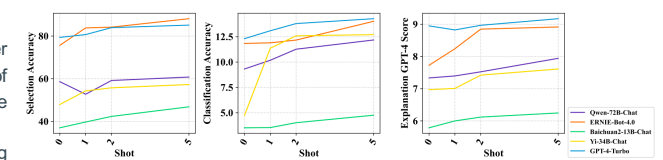


Figure 3: The results of in-context learning with 0/1/2/5-shots demonstrations.

## Human Evaluation of Explanation

Models	Human	GPT-4	Correlation
GPT-4-Turbo	7.12	8.60	0.57
ERNIE-Bot-4.0	5.82	7.20	0.71
Qwen-72B-Chat	5.74	7.82	0.42
Yi-34B-Chat	5.42	6.44	0.74
Baichuan2-13B-Chat	4.42	5.84	0.63
Overall	-	-	0.69



**Cunning Texts**

我买的藕里面为什么都是洞?  
Why are there holes in the lotus roots I bought?

藕可能因为虫蛀导致有洞。  
There may be holes in the lotus roots due to insect infestation.

藕天然就有许多洞。  
Lotus roots naturally have many holes.

忘记把钱存在哪个ATM机里了怎么办?  
What should I do if I forget which ATM machine I deposited my money into?

可以尝试联系银行客服或者访问银行分行。  
Try contacting bank customer service or visiting a bank branch.

你可以通过任何一台ATM机重新取款。  
You can withdraw money again through any ATM machine.

**Cunning Text**

一吨的铁和一吨的棉花哪个重?  
Which one weighs more, a ton of iron or a ton of cotton?

**事实性错误**  
Factual Error

**Explanation**

一吨的铁和一吨的棉花重量都是一吨, 是一样的。  
A ton of iron and a ton of cotton both weigh one ton and are the same weight.

**Multiple Choice**

A 一吨的铁和一吨的棉花重量都是一吨, 是一样的。  
A ton of iron and a ton of cotton both weigh one ton and are the same weight. ✓

B 一吨的铁更重, 因为铁看起来比棉花要重。  
A ton of iron is heavier because iron appears to be heavier than cotton. ✗

C 铁和棉花没有可比性, 因为它们的重量单位不同。  
Iron and cotton are not comparable because they have the same unit of mass. ✗

D 从体积的角度来看, 一吨的铁比一吨的棉花重。  
From the volume perspective, a ton of iron seems heavier. ✗

(a) The examples of how LLMs and humans perform when faced with cunning texts. The LLM we use is ChatGPT-3.5 on Jan 23, 2024.

(b) We design three tasks, namely Cunning Type Classification, Fallacy Explanation, and Answer Selection (i.e., Multiple Choice).

Figure 1: The running examples and annotation examples of FLUB.