# Are Large Language Models Really Good Logical Reasoners? A Comprehensive Evaluation and Beyond

Fangzhi Xu*, Qika Lin*, Jiawei Han, Tianzhe Zhao, Jun Liu, *Senior Member, IEEE,*
Erik Cambria, *Fellow, IEEE*

**Abstract**—Logical reasoning consistently plays a fundamental and significant role in the domains of knowledge engineering and artificial intelligence. Recently, Large Language Models (LLMs) have emerged as a noteworthy innovation in natural language processing (NLP), exhibiting impressive achievements across various classic NLP tasks. However, the question of whether LLMs can effectively address the task of logical reasoning, which requires gradual cognitive inference similar to human intelligence, remains unanswered. To this end, we aim to bridge this gap and provide comprehensive evaluations in this paper. Firstly, to offer systematic evaluations, we select fifteen typical logical reasoning datasets and organize them into deductive, inductive, abductive and mixed-form reasoning settings. Considering the comprehensiveness of evaluations, we include three representative LLMs (i.e., text-davinci-003, ChatGPT and BARD) and evaluate them on all selected datasets under zero-shot, one-shot and three-shot settings. Secondly, different from previous evaluations relying only on simple metrics (e.g., *accuracy*), we propose fine-level evaluations from objective and subjective manners, covering both answers and explanations. These evaluations include *answer correctness*, *explain correctness*, *explain completeness* and *explain redundancy*. Additionally, to uncover the logical flaws of LLMs, problematic cases will be attributed to five error types from two dimensions, i.e., *evidence selection process* and *reasoning process*. The former one includes *wrong selection* and *hallucination*, while the latter one includes *no reasoning*, *perspective mistake* and *process mistake*. Thirdly, to avoid the influences of knowledge bias and concentrate purely on benchmarking the logical reasoning capability of LLMs, we propose a new dataset with neutral content. It contains 3,000 samples and covers deductive, inductive and abductive settings. Based on the in-depth evaluations, this paper finally forms a general evaluation scheme of logical reasoning capability from six dimensions (i.e., *Correct*, *Rigorous*, *Self-aware*, *Active*, *Oriented* and *No hallucination*). It reflects the pros and cons of LLMs and gives guiding directions for future works.

**Index Terms**—Logical reasoning, large language model, deductive reasoning, inductive reasoning, abductive reasoning.

---◆---

## 1 INTRODUCTION

As a fundamental and significant topic in the domains of knowledge engineering and artificial intelligence, logical reasoning has consistently remained a subject of intense research interest [1], [2], [3]. Through the integration of logical reasoning, a wide range of intelligent applications can be developed (e.g., recommendation systems [4], relation prediction [5] and question generation [6]), which not only offer powerful capabilities but also ensure natural interpretability. Nevertheless, the development of efficient and robust logical reasoning systems continues to be a challenging academic pursuit, primarily attributed to the complexities associated with handling intricate semantic and syntactic structures, managing symbolic background knowledge, executing advanced high-level abstraction and

inference processes, as well as navigating through substantial uncertainty and ambiguous information [7], [8].

With the aid of large-scale pre-training, instruction fine-tuning, and human feedback reinforcement learning strategies [9], Large Language Models (LLMs) have made unparalleled strides in the artificial intelligence community, particularly in the field of natural language processing (NLP). They have achieved remarkable performance in numerous traditional NLP tasks [10], [11], such as question answering, information retrieval, machine translation, and affective computing. Under this circumstance, researchers have begun to question the efficacy of LLMs in addressing complex logical reasoning tasks and their ability to apply this knowledge to downstream intelligent applications. As such, a natural question arises: *are large language models really good logical reasoners*?

There are already several studies to evaluate the capability of LLMs from various reasoning perspectives, e.g., multilingual reasoning [12], commonsense reasoning [13], and mathematical reasoning [14]. These efforts in evaluating the specific capabilities of LLMs are meaningful and can benefit future directions. Nevertheless, all of them are confronted with one or more of the following defects: (1) There are no comprehensive evaluations, which are hampered by the issue of lacking systematic category, limited LLMs for comparison, and limited data samples for evaluation; (2)

- *Fangzhi Xu, Qika Lin, Jiawei Han and Tianzhe Zhao are with the School of Computer Science and Technology, Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China (email: {Leo981106, tara1208260223}@stu.xjtu.edu.cn, {qikalin, ztz8758}@foxmail.com).*
- *Jun Liu is with the Shaanxi Provincial Key Laboratory of Big Data Knowledge Engineering, and National Engineering Lab for Big Data Analytics, Xi'an, Shaanxi 710049, China (e-mail: liukeen@xjtu.edu.cn).*
- *Erik Cambria is with the School of Computer Science and Engineering, Nanyang Technological University (e-mail: cambria@ntu.edu.sg).*
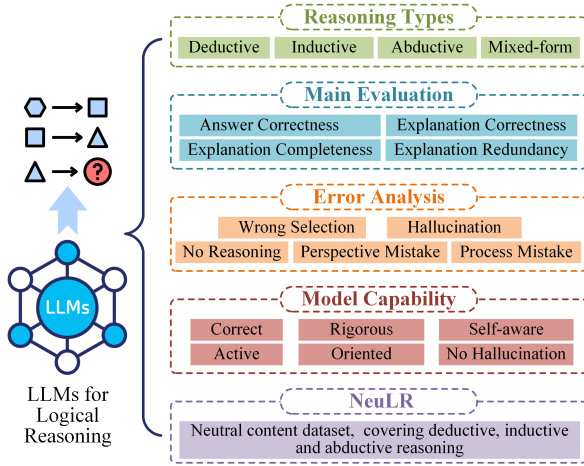- *Fangzhi Xu and Qika Lin contribute equally to this paper.*

Fig. 1: The overall architecture of the evaluation.

The majority of evaluation works purely report the accuracy of answers, which limits fine-level analysis on the reasoning process and fails to explore the causes of mistakes; (3) Current logical reasoning benchmarks may fail to purely evaluate logical reasoning ability since the reasoning of LLMs can be affected by the content; (4) There lack a complete evaluation system or well-defined dimensions to comprehensively conclude the logical reasoning capability of LLMs. Therefore, our work aims to fill these gaps and provide a comprehensive evaluation, where the overall architecture encompasses the five main points as illustrated in Fig. 1: category with reasoning settings, main metrics for evaluation, perspectives of error analysis, dimensions for measuring logical reasoning capability, and new benchmark NeuLR with neutral content. In detail, we address the above limitations from the following aspects.

Firstly, our work starts from systematic views and provides comprehensive evaluations. According to the classical definition [15], logical reasoning can be mainly categorized into three fundamental types, i.e., deductive, inductive and abductive reasoning. They together form a complete chain of reasoning, thus it is meaningful to evaluate LLMs from these views. Based on it, our comprehensive evaluation is reflected in three orthogonal dimensions. For the reasoning type view, all the evaluated datasets are categorized into four reasoning settings, i.e., deductive, inductive, abductive and mixed-forms. The former three involve the independent reasoning manner. Considering some recent efforts have been made on proposing challenging settings with mixed reasoning manners, we introduce an additional category of *mixed-form* for evaluations. For the dataset view, we include fifteen typical logical reasoning datasets according to the above categories. For the model view, we evaluate three representative LLMs, i.e., text-davinci-003 [16], ChatGPT [17] and BARD [18].

Secondly, our work fills the blank in fine-level evaluations of logical reasoning tasks. Current benchmarks only rely on a few objective metrics (e.g., accuracy) to measure the model capability. It may not be sufficient in the case of generative LLMs, since the role of LLMs is not only limited to correctly answer questions but also serves as practical tools, which are required to provide reasoning chains or explanations. Previous works [19], [20] conduct extensive experiments on popular NLP datasets, but they purely report the performance results. Since some LLMs (e.g., Chat-GPT) function as the interactive tools for human use, it is necessary to introduce subjective metrics to do fine-grained evaluations. In this paper, we employ four dimensions of metrics, covering *answer correctness*, *explanation correctness*, *explanation completeness* and *explanation redundancy*. It can provide more meaningful and complete evaluations from both objective and subjective views. Considering problematic cases (i.e., wrong answer or wrong explanation) can reflect obvious logical flaws of LLMs, we further attribute them to several error types from two dimensions of *evidence selection process* and *reasoning process* and give in-depth analysis.

Thirdly, our work focuses on the issue of content neutrality and provides new solutions. The current benchmark for evaluating logical reasoning ability is strongly coupled with text comprehension; in other words, the rule reasoning process of LLMs may be affected by the content in the inputs, which limits the test of real logical reasoning ability. Also, LLMs are highly powerful due to their massive training data, which may overlap with popular benchmarks. As a result, testing LLMs on these benchmarks may not be entirely fair, as it can only demonstrate the fitting ability of LLMs rather than their real logical reasoning capability. Therefore, language models may be trained to learn a biased pattern from text, rather than really capture the logical reasoning capability. Some previous works [21], [22] propose to establish complete benchmarks for LLMs. But few works focus on logical reasoning and fail to attend to the content-neutral problem. To narrow this gap, we propose a new dataset named NeuLR, which contains 3,000 content-neutral samples and covers the deductive, inductive and abductive reasoning types. It is expected to offer a novel perspective for benchmarking the logical reasoning ability of LLMs.

Finally, we conclude the extensive performance results of LLMs and form an evaluation scheme with six key properties, i.e., *Correct*, *Rigorous*, *Self-aware*, *Active*, *Oriented* and *No Hallucination*. Among them, *Correct* purely measures the accuracy of the answer. *Rigorous* measures whether LLMs give both correct answers and complete and correct explanations. *Self-aware* is reflected by the redundancy of the generated content. *Active* is measured by the proportion of reasoning. *Oriented* illustrates whether LLMs can reason from the right perspectives. *No hallucination* measures whether LLMs are more prone to produce hallucination. The above dimensions can all be quantified from existing evaluation experiments. For deductive, inductive, abductive and mixed reasoning settings respectively, we obtain the ability maps based on the six properties for each LLM. It is meaningful to identify the strengths and weaknesses of LLMs under the four reasoning settings, thus guiding future directions.

The main contributions of the paper are listed as follows:

(1) In view of the great success of LLMs in massive NLP tasks, our work is targeted at answering *are LLMs really good logical reasoners?*. In this paper, we provide a comprehensive evaluation and give potential directions for future researches.

(2) For a comprehensive evaluation of logical reasoning,

this paper classifies datasets into four reasoning manners, i.e., deductive, inductive, abductive and mixed-form. In light of the insufficient evaluations of previous works, we include fifteen typical logical reasoning datasets and evaluate three representative LLMs (i.e., text-davinci-003, ChatGPT and BARD) under both zero- and few-shot settings.

(3) Considering the drawbacks in current objective metrics, this paper gives fine-level evaluations including four dimensions i.e., answer correctness, explanation correctness, explanation completeness and explanation redundancy. To explore the value of failure cases, we attribute them to several error types and explore the logical flaws of LLMs.

(4) To provide fair evaluations with neutral content and decouple logical reasoning from text understanding, this paper proposes a new dataset named NeuLR[1]. It contains 3,000 content-neutral samples and covers deductive, inductive and abductive reasoning manners.

(5) In view of the evaluation results, this paper forms a general evaluation scheme for the logical reasoning capability of LLMs for the first time, which concludes six key properties, i.e, *Correct*, *Rigorous*, *Self-aware*, *Active*, *Oriented* and *No hallucination*. Furthermore, we derive the ability maps for each LLM under four reasoning settings respectively and propose future directions.

## 2 PRELIMINARY

Logical reasoning aims to generate logical implications that contain new facts using one-step or multi-step inference based on given premises [23], [24], [25], i.e., *premise⇒conclusion*. Elements of logical reasoning typically include knowledge facts or premises and logical rules, for example:

- rule: Children of eight years old are all in primary school.
- fact1/premise1: Jordan is a child of eight years old.
- fact2/premise2: Jordan is in primary school.

According to the reasoning classification system of classical logic [15], there are three major types of logical reasoning: deductive, inductive and abductive. Drawing upon the aforementioned rule and facts, we can depict the objective of these three types of reasoning as forecasting the third item based on the two provided ones.

**Deductive Reasoning**. Deductive reasoning is the psychological process of drawing deductive inferences that start from the given premises and reason with logical rules or commonsense to obtain *certain* conclusions [26], [27]. It can be *premise1+rule⇒premise2*. Fig. 2 presents an example of deductive reasoning. Its reasoning progress generates specific knowledge facts from general counterparts, e.g., *premise2* and *rule* are specific and general knowledge, respectively. Therefore, deductive reasoning is actually a top-down way.

**Inductive Reasoning**. Distinct from deductive reasoning, inductive reasoning derives general principles from a body of observations which means making broad generalizations based on specific observations [28], [29]. For example, an example of inductive reasoning can be *premise1+premise2⇒rule*, concluding generalized knowledge

Fig. 2: An example of using an LLM to answer reasoning questions. The red words represent the generated results, while the blue ones represent the generated explanations.

*rule* that is independent with specific item *Jordan*. Generally, the truth of the conclusion of an inductive argument is *probable* rather than *certain* in inductive reasoning. Thus, inductive reasoning is a bottom-up approach and is contrasted with deductive reasoning.

**Abductive Reasoning**. Formally, abductive reasoning is similar to deductive reasoning which seeks conclusions from a set of observations. But differently, its target is to generate the simplest and most likely explanation for the given observations [30], [31]. So the result is *probable* like in inductive reasoning. An example of abductive reasoning can be *premise2+rule⇒premise1*, which means *premise1* is the most likely cause of *premise2*.

These three types of reasoning encompass the fundamental patterns of logic. Nevertheless, numerous real-life reasoning scenarios may require multiple inference steps and the incorporation of at least two of these three types. In this paper, they can be considered as a more intricate type of reasoning, namely, *mixed-form*.

## 3 EVALUATION DETAILS

This section presents the comprehensive experimental settings, including the models, datasets and metrics. More detailed prompts for zero- or few-shot are shown in Table 3 of the Appendix [2].

### 3.1 Evaluated Models

Because of the rapid emergence of LLMs, it is not realistic to include all LLMs in this paper. Thus, we select three representative ones for evaluation, which are text-davinci-003 [3], ChatGPT [4] and BARD [5]. The details of these three models are listed in Table 1 of the Appendix.

Among them, text-davinci-003 is the earliest LLM released by OpenAI, which is expected to undertake any

language task. For ChatGPT, we utilize the version of *GPT-3.5-turbo* for evaluation, which is the most capable and cost-effective version in the GPT-3.5 family. BARD flatform is based on one of the latest LLM, i.e., PaLM 2 [18], which is updated and released by Google in May 2023. Also, it is relatively larger in size compared with the GPT-3.5 family.

## 3.2 Evaluated Datasets

According to the previous discussion, the evaluation is conducted systematically from deductive, inductive, abductive and mixed-form views. Therefore, this paper selects fifteen popular datasets in logical reasoning and divides them into the above four folds. Table 2 in the Appendix presents detailed information of these datasets. The selected datasets contain both generation and classification ones and there exist diverse forms of tasks, which illustrate the comprehensiveness of our evaluation. Diverging from prior works that only employ a limited number of samples, this paper significantly expands the evaluation size. Since ChatGPT is one of the most popular LLM for the public, we give much focus on it. For parts of the datasets, we keep all the test examples for ChatGPT evaluation (i.e., EntailmentBank [32], FOLIO [33], Leap-Of-Thought [34], CLUTRR [35], ReClor [36], LogiQA [37], LogiQA 2.0 [38] and LogiQA2NLI [38]), while other large datasets (i.e., bAbI-15 [39], RuleTaker [40], bAbI-16 [39], $\alpha$-NLI [41], $\alpha$-NLG [41], AbductiveRules [42] and D*-Ab [43]) are sampled to 1,000 examples. As for the evaluation of text-davinci-003 and BARD, we sample 100 test examples for each dataset.

## 3.3 Selected Metrics

The majority of prior evaluation studies solely report the accuracy metric, which may not be comprehensive enough for evaluating the effectiveness of LLMs. Consequently, we assert that a more nuanced assessment is necessary. To this end, we propose to evaluate LLMs from both objective and subjective perspectives. Specifically, we introduce four evaluation metrics to reflect the intermediate reasoning process of LLMs: *answer correctness*, *explanation correctness*, *explanation completeness* and *explanation redundancy*.

- **Answer Correctness**. This metric ascertains the degree of consistency between the generated answer and the true label. In the context of generation tasks, it necessitates that the meanings of the two outputs match, instead of their corresponding tokens.
- **Explanation Correctness**. It indicates whether the generated explanation is logically correct to reason towards the true answer. It is a subjective view to determine whether the reasoning process of machines is in line with that of humans.
- **Explanation Completeness**. It means that in the reasoning process, the correct answer can be inferred through the selected known facts and the generated intermediate facts. This does not necessarily cause answer correctness or explanation correctness.
- **Explanation Redundancy**. It implies that the chosen established facts and the generated intermediate facts are in excess of the requisite facts needed to ascertain an accurate conclusion. Consequently, this

introduces superfluous and extraneous facts into the reasoning process.

Notably, the metric values of explanation correctness, completeness, and redundancy are independent.

To identify prevalent logical flaws in LLMs, we establish error types to categorize problematic cases. This paper classifies errors along two primary dimensions: (1) *evidence selection process* and (2) *reasoning process*. The first dimension centers on evaluating the evidence selected by LLMs, whereas the second dimension emphasizes the logical reasoning process using the selected evidence. Detailedly, *evidence selection process* category can be further divided into *wrong selection* and *hallucination*. The former denotes that LLMs select the wrong facts or ignore the necessary facts from the beginning of the reasoning. The latter denotes that LLMs select the evidence which contradicts the given context or can not be verified by the context. *Reasoning process* category can be further divided into *no reasoning*, *perspective mistake* and *process mistake*. The first error type signifies instances where LLMs fail to conduct reasoning, instead merely listing the given facts and the final answer. The second denotes LLMs starting from an irrelevant point or focusing on an improper perspective for the correct answer. The last refers to LLMs commencing from a proper viewpoint, but making mistakes during the reasoning process.

## 4 OVERALL EXPERIMENTS

In this section, we conduct evaluation experiments on three LLMs, i.e., text-davinci-003, ChatGPT and BARD under zero-shot, one-shot and three-shot settings respectively. Table 1 presents the overall answer correctness of these three LLMs on fifteen logical reasoning datasets. Generally, LLMs' performances on logical reasoning tasks still have significant room for improvement in comparison to the state-of-the-art (SOTA) metric. Most of the results fall short of those achieved by smaller-sized SOTA models. We provide a detailed analysis of the results from the following perspectives.

Firstly, we conduct an analysis of LLMs' performances across four reasoning manners, solely focusing on the zero-shot results to facilitate a clear comparison. Furthermore, we introduce the relative performance metric (i.e., LLM accuracy/SOTA) to reflect the relative capability of LLMs in comparison to SOTA performances in Fig. 3. We also calculate the weighted results of four reasoning manners in Fig. 4a. From the results, ChatGPT performs worse in deductive and inductive settings compared with text-davinci-003 and BARD. In the abductive setting, three LLMs show comparable performances and BARD wins with slight advantages. In the mixed-form setting, ChatGPT performs better and BARD ranks second. Overall, BARD shows consistent superiority among deductive, inductive and abductive settings, while text-davinci-003 also does relatively well. It seems that ChatGPT struggles in the three settings, but is better at mixed-form reasoning. Also, we compare the LLM performances between deductive, inductive and abductive settings. LLMs do best in deductive setting, while they mostly struggle in inductive setting. We argue that deductive and abductive reasoning align with typical NLP scenarios, where LLMs have to provide missing facts. Conversely, inductive reasoning necessitates extracting high-level rules

TABLE 1: Overall results of LLMs' answer correctness across the zero-shot, one-shot and three-shot logical reasoning settings. The notations *De.*, *In.*, *Ab.* and *Mix* correspond to deductive, inductive, abductive and mixed-form reasoning, respectively (as in the following tables and figures). *Gen.* indicates whether the task is a generation one. The percentage signs (%) of performance values are omitted for simplicity in the paper.

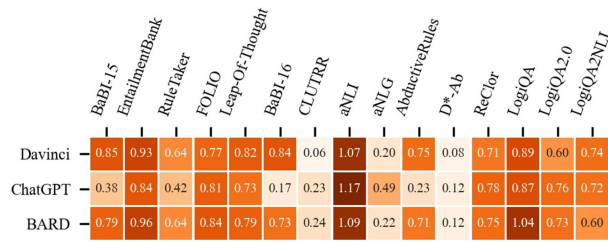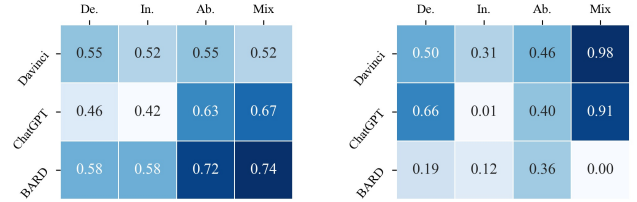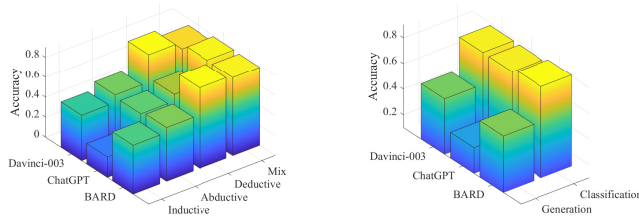| | Dataset | Gen. | text-davinci-003 | | | ChatGPT | | | BARD | | | SOTA |
| | | | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **De.** | **bAbI-15** | ✓ | 85.00 | 76.00 | 75.00 | 38.40 | 46.40 | 39.70 | 79.00 | 80.00 | 88.00 | 100 [39] |
| | **EntailmentBank** | ✓ | 93.00 | 88.00 | 89.00 | 83.82 | 82.06 | 77.94 | 96.00 | 97.00 | 97.00 | 100 [32] |
| | **RuleTaker** | | 64.00 | 60.00 | 62.00 | 42.00 | 38.00 | 40.20 | 64.00 | 57.00 | 70.00 | ≈100 [40] |
| | **FOLIO** | | 48.00 | 53.00 | 52.00 | 50.00 | 50.98 | 54.41 | 52.00 | 43.00 | 49.00 | 62.11 [33] |
| | **Leap-Of-Thought** | | 82.00 | 90.00 | 87.00 | 72.61 | 74.01 | 61.21 | 79.00 | 72.00 | 79.00 | 99.7 [34] |
| **In.** | **bAbI-16** | ✓ | 84.00 | 81.00 | 74.00 | 17.10 | 24.70 | 12.90 | 73.00 | 44.00 | 52.00 | 100 [39] |
| | **CLUTRR** | ✓ | 6.00 | 23.00 | 20.00 | 21.99 | 19.55 | 12.83 | 23.00 | 26.00 | 24.00 | 95.0 [44] |
| **Ab.** | **α-NLI** | | 74.00 | 70.00 | 74.00 | 80.90 | 80.00 | 79.10 | 75.00 | 74.00 | 77.00 | 68.90 [41] |
| | **α-NLG** | ✓ | 9.00 | 10.00 | 12.00 | 21.90 | 23.40 | 25.90 | 10.00 | 12.00 | 15.00 | 45.00 [41] |
| | **AbductiveRules** | ✓ | 75.00 | 42.00 | 35.00 | 23.30 | 35.10 | 29.80 | 71.00 | 49.00 | 22.00 | 100 [42] |
| | **D*-Ab** | ✓ | 8.00 | 21.00 | 23.00 | 11.60 | 2.50 | 1.80 | 11.00 | 0.00 | 0.00 | ≥95 [43] |
| **Mix** | **ReClor** | | 53.00 | 53.00 | 55.00 | 58.80 | 56.00 | 58.80 | 56.00 | 55.00 | 56.00 | 75.00 [45] |
| | **LogiQA** | | 41.00 | 35.00 | 39.00 | 40.25 | 39.48 | 40.86 | 48.00 | 46.00 | 47.00 | 46.10 [45] |
| | **LogiQA 2.0** | | 43.00 | 42.00 | 41.00 | 54.60 | 50.80 | 54.80 | 53.00 | 46.00 | 47.00 | 72.25 [38] |
| | **LogiQA2NLI** | | 59.00 | 55.00 | 58.00 | 57.83 | 53.83 | 57.00 | 48.00 | 50.00 | 47.00 | ≈80 |

Fig. 3: LLM performances on different datasets.

(a) Rigorous evaluation.     (b) Self-aware evaluation.

Fig. 5: Heatmap visualization of rigor and self-awareness.

(a) Different reasoning types.     (b) Generation/Classification.

Fig. 4: Visualization on the metric of answer correctness.

or knowledge from the given facts, which is more intricate and may not be readily available in the training corpus.

Secondly, we conduct a detailed analysis of LLMs' performances from the generation and classification perspectives in Fig. 4b. In general, classification scenarios tend to yield better performance than generation counterparts. Notably, ChatGPT exhibits particularly poor results in generation tasks, such as bAbI-15, bAbI-16, CLUTRR, AbductiveRules and D*-Ab. This observation may result from the fact that ChatGPT is designed to improve chatting capability rather than complex reasoning, which can lead to performance degradation in pure generative logic reasoning scenarios.

Thirdly, few-shot in-context learning (ICL) [46] does not necessarily bring improvements in logical reasoning tasks.

It is quite inconsistent with the cases in other non-reasoning NLP tasks, such as topic classification and sentiment analysis [47]. We count the cases where LLMs can continuously obtain the performance gains from few-shot ICL (i.e., 0-shot < 1-shot < 3-shot). For text-davinci-003, only two (out of four) abductive datasets continuously benefit from the few-shot ICL. ChatGPT witnesses performance improvements only in one (out of five) deductive dataset and one (out of four) abductive dataset. For BARD, few-shot ICL helps two (out of five) deductive datasets and one (out of four) abductive datasets. Remarkably, few-shot ICL fails to provide consistent benefits for LLMs under inductive reasoning and mixed-form reasoning manners. We argue that inductive and mixed-form settings require more complex and high-order reasoning ability, which may be difficult to learn with few samples and the ICL samples may cause noises. But the task form of deductive and abductive reasoning is easy to follow, which provides the application potential for few-shot ICL.

## 5 FINE-LEVEL EVALUATIONS

This section presents a detailed evaluation of LLMs from various perspectives. Firstly, we concentrate on our proposed four metrics and offer comprehensive analyses from diverse dimensions. Secondly, we select problematic cases and attribute the reasoning errors. Thirdly, we conduct

TABLE 2: Evaluations on whether LLMs are rigorous reasoners. For each dataset, the first row of results represents the performances when LLMs give the correct answer, correct explanation as well as complete explanation simultaneously. The values in the subscripts denote the drops compared with only distinguishing the answer correctness. The second row of results represents the performances when LLMs give both correct answers and correct explanations, regardless of the explanation completeness. The third row represents the cases when LLMs give correct answers and list complete explanations, regardless of the explanation correctness.

| | Dataset | text-davinci-003 | | | ChatGPT | | | BARD | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot |
| **Deductive** | **bAbI-15** | $53.00_{32.00\downarrow}$ | $60.00_{16.00\downarrow}$ | $64.00_{11.00\downarrow}$ | $25.50_{12.90\downarrow}$ | $12.10_{34.30\downarrow}$ | $14.10_{25.60\downarrow}$ | $45.00_{34.00\downarrow}$ | $25.00_{55.00\downarrow}$ | $47.00_{41.00\downarrow}$ |
| | | 61.00 | 66.00 | 68.00 | 32.10 | 12.90 | 16.40 | 77.00 | 74.00 | 85.00 |
| | | 56.00 | 60.00 | 64.00 | 27.90 | 18.00 | 17.60 | 45.00 | 25.00 | 47.00 |
| | **EntailmentBank** | $29.00_{64.00\downarrow}$ | $37.00_{51.00\downarrow}$ | $30.00_{59.00\downarrow}$ | $25.88_{57.94\downarrow}$ | $20.00_{62.06\downarrow}$ | $10.59_{67.35\downarrow}$ | $26.00_{38.00\downarrow}$ | $25.00_{32.00\downarrow}$ | $33.00_{37.00\downarrow}$ |
| | | 29.00 | 37.00 | 30.00 | 25.88 | 20.00 | 10.59 | 54.00 | 66.00 | 71.00 |
| | | 72.00 | 73.00 | 75.00 | 62.06 | 57.65 | 31.76 | 94.00 | 96.00 | 97.00 |
| | **RuleTaker** | $35.30_{6.70\downarrow}$ | $22.50_{15.50\downarrow}$ | $24.80_{15.40\downarrow}$ | $25.88_{57.94\downarrow}$ | $20.00_{62.06\downarrow}$ | $10.59_{67.35\downarrow}$ | $54.00_{42.00\downarrow}$ | $66.00_{31.00\downarrow}$ | $71.00_{26.00\downarrow}$ |
| | | 36.20 | 24.00 | 26.00 | 25.88 | 20.00 | 10.59 | 26.00 | 28.00 | 33.00 |
| | | 36.00 | 23.20 | 26.00 | 62.06 | 57.65 | 31.76 | 26.00 | 25.00 | 34.00 |
| | **FOLIO** | $27.00_{21.00\downarrow}$ | $25.00_{28.00\downarrow}$ | $25.00_{27.00\downarrow}$ | $28.92_{21.08\downarrow}$ | $27.94_{23.04\downarrow}$ | $27.94_{26.47\downarrow}$ | $21.00_{31.00\downarrow}$ | $19.00_{24.00\downarrow}$ | $20.00_{29.00\downarrow}$ |
| | | 28.00 | 25.00 | 26.00 | 28.92 | 27.94 | 27.94 | 23.00 | 19.00 | 22.00 |
| | | 27.00 | 28.00 | 27.00 | 33.33 | 32.35 | 35.29 | 25.00 | 22.00 | 23.00 |
| | **Leap-of-Thought** | $29.00_{53.00\downarrow}$ | $43.00_{47.00\downarrow}$ | $38.00_{49.00\downarrow}$ | $70.60_{2.02\downarrow}$ | $24.36_{49.65\downarrow}$ | $28.86_{32.35\downarrow}$ | $76.00_{3.00\downarrow}$ | $69.00_{3.00\downarrow}$ | $77.00_{2.00\downarrow}$ |
| | | 63.00 | 63.00 | 58.00 | 71.22 | 24.83 | 29.79 | 76.00 | 69.00 | 77.00 |
| | | 32.00 | 46.00 | 42.00 | 70.99 | 48.33 | 40.88 | 79.00 | 71.00 | 7.00 |
| **Inductive** | **bAbI-16** | $59.00_{25.00\downarrow}$ | $35.00_{46.00\downarrow}$ | $23.00_{51.00\downarrow}$ | $8.30_{8.80\downarrow}$ | $8.20_{16.50\downarrow}$ | $2.60_{10.30\downarrow}$ | $24.00_{49.00\downarrow}$ | $15.00_{29.00\downarrow}$ | $16.00_{36.00\downarrow}$ |
| | | 67.00 | 50.00 | 37.00 | 10.20 | 8.40 | 3.10 | 58.00 | 32.00 | 24.00 |
| | | 65.00 | 44.00 | 35.00 | 10.00 | 9.40 | 3.50 | 32.00 | 22.00 | 24.00 |
| | **CLUTRR** | $2.00_{4.00\downarrow}$ | $7.00_{16.00\downarrow}$ | $6.00_{14.00\downarrow}$ | $7.85_{14.14\downarrow}$ | $4.62_{14.92\downarrow}$ | $2.71_{10.12\downarrow}$ | $19.00_{4.00\downarrow}$ | $25.00_{1.00\downarrow}$ | $23.00_{1.00\downarrow}$ |
| | | 2.00 | 7.00 | 6.00 | 7.94 | 4.62 | 2.71 | 19.00 | 25.00 | 23.00 |
| | | 6.00 | 18.00 | 18.00 | 10.56 | 6.72 | 3.75 | 21.00 | 26.00 | 23.00 |
| **Abductive** | **$\alpha$-NLI** | $68.00_{6.00\downarrow}$ | $69.00_{1.00\downarrow}$ | $68.00_{6.00\downarrow}$ | $77.50_{3.40\downarrow}$ | $64.50_{15.50\downarrow}$ | $58.40_{20.70\downarrow}$ | $75.00_{0.00-}$ | $71.00_{3.00\downarrow}$ | $77.00_{0.00-}$ |
| | | 70.00 | 69.00 | 68.00 | 78.20 | 66.70 | 60.90 | 75.00 | 71.00 | 77.00 |
| | | 68.00 | 69.00 | 68.00 | 77.60 | 64.50 | 58.40 | 75.00 | 71.00 | 77.00 |
| | **$\alpha$-NLG** | $1.00_{8.00\downarrow}$ | $0.00_{10.00\downarrow}$ | $2.00_{10.00\downarrow}$ | $15.30_{6.60\downarrow}$ | $16.00_{7.40\downarrow}$ | $10.30_{15.60\downarrow}$ | $7.00_{3.00\downarrow}$ | $8.00_{4.00\downarrow}$ | $9.00_{6.00\downarrow}$ |
| | | 1.00 | 0.00 | 2.00 | 15.30 | 16.00 | 10.40 | 7.00 | 9.00 | 10.00 |
| | | 8.00 | 8.00 | 9.00 | 20.90 | 22.40 | 21.40 | 10.00 | 9.00 | 14.00 |
| | **AbductiveRules** | $50.00_{25.00\downarrow}$ | $5.00_{37.00\downarrow}$ | $0.00_{35.00\downarrow}$ | $12.00_{11.30\downarrow}$ | $18.30_{16.80\downarrow}$ | $5.70_{24.10\downarrow}$ | $57.00_{14.00\downarrow}$ | $30.00_{19.00\downarrow}$ | $10.00_{12.00\downarrow}$ |
| | | 75.00 | 10.00 | 0.00 | 20.50 | 29.60 | 9.50 | 65.00 | 40.00 | 18.00 |
| | | 50.00 | 5.00 | 0.00 | 13.00 | 19.90 | 5.80 | 57.00 | 30.00 | 10.00 |
| | **D*-Ab** | $4.00_{4.00\downarrow}$ | $5.00_{16.00\downarrow}$ | $4.00_{19.00\downarrow}$ | $4.10_{7.50\downarrow}$ | $1.20_{1.30\downarrow}$ | $1.10_{0.70\downarrow}$ | $4.00_{7.00\downarrow}$ | $0.00_{0.00-}$ | $0.00_{0.00-}$ |
| | | 7.00 | 7.00 | 7.00 | 5.50 | 1.50 | 1.20 | 4.00 | 0.00 | 0.00 |
| | | 4.00 | 5.00 | 4.00 | 4.50 | 1.20 | 1.10 | 5.00 | 0.00 | 0.00 |
| **Mixed-form** | **ReClor** | $5.00_{48.00\downarrow}$ | $0.00_{53.00\downarrow}$ | $0.00_{55.00\downarrow}$ | $28.60_{30.20\downarrow}$ | $25.20_{30.80\downarrow}$ | $29.80_{29.00\downarrow}$ | $38.00_{18.00\downarrow}$ | $34.00_{21.00\downarrow}$ | $33.00_{23.00\downarrow}$ |
| | | 5.00 | 0.00 | 0.00 | 32.40 | 28.00 | 32.20 | 48.00 | 46.00 | 42.00 |
| | | 42.00 | 46.00 | 42.00 | 50.00 | 32.80 | 46.00 | 38.00 | 36.00 | 38.00 |
| | **LogiQA** | $5.00_{36.00\downarrow}$ | $2.00_{33.00\downarrow}$ | $1.00_{38.00\downarrow}$ | $25.96_{14.29\downarrow}$ | $21.35_{18.13\downarrow}$ | $20.43_{20.43\downarrow}$ | $29.00_{19.00\downarrow}$ | $26.00_{20.00\downarrow}$ | $19.00_{28.00\downarrow}$ |
| | | 6.00 | 2.00 | 1.00 | 28.42 | 23.20 | 23.81 | 34.00 | 27.00 | 23.00 |
| | | 23.00 | 26.00 | 27.00 | 29.19 | 27.50 | 26.42 | 30.00 | 35.00 | 27.00 |
| | **LogiQA2.0** | $38.00_{5.00\downarrow}$ | $32.00_{10.00\downarrow}$ | $28.00_{13.00\downarrow}$ | $43.40_{11.20\downarrow}$ | $34.60_{16.20\downarrow}$ | $38.80_{16.00\downarrow}$ | $47.00_{6.00\downarrow}$ | $39.00_{7.00\downarrow}$ | $44.00_{3.00\downarrow}$ |
| | | 39.00 | 32.00 | 28.00 | 43.80 | 35.20 | 38.80 | 47.00 | 39.00 | 44.00 |
| | | 39.00 | 32.00 | 30.00 | 44.20 | 36.80 | 40.60 | 47.00 | 39.00 | 44.00 |
| | **LogiQA2NLI** | $57.00_{2.00\downarrow}$ | $51.00_{4.00\downarrow}$ | $56.00_{2.00\downarrow}$ | $43.17_{14.67\downarrow}$ | $36.33_{17.50\downarrow}$ | $36.50_{20.50\downarrow}$ | $38.00_{10.00\downarrow}$ | $41.00_{9.00\downarrow}$ | $37.00_{10.00\downarrow}$ |
| | | 57.00 | 51.00 | 56.00 | 43.50 | 36.50 | 36.83 | 40.00 | 43.00 | 40.00 |
| | | 57.00 | 51.00 | 56.00 | 53.50 | 49.00 | 50.16 | 38.00 | 41.00 | 37.00 |

additional analyses on specific datasets and present some phenomena observed in different reasoning manners.

## 5.1 Are LLMs Rigorous Logical Reasoning?

While LLMs may produce correct answers in some cases, it is unclear whether they perform the correct logical reasoning or simply arrive at the right answer by chance. Therefore, we delve deeper into the reasoning process beyond the output answers. We view cases where LLMs provide a correct answer along with a correct and complete explanation as *Rigorous*. The detailed results are shown in Table 2. Compared with the simple judgment of answer correctness, all selected LLMs present obvious performance drops. To simplify the analysis, we only take the zero-shot setting into consideration. In Fig. 5a, we calculate the ratio of rigorous performance and answer accuracy. The higher values (darker colors) mean better performance in rigorous reasoning. According to the results, BARD shows the best capability in rigorous reasoning, consistently under four reasoning manners. Meanwhile, ChatGPT still struggles in deductive and inductive settings, while text-davinci-003 comes last in both abductive and mixed-form manners.

TABLE 3: Evaluation results on the metric of explanation redundancy.

| | Dataset | Gen. | text-davinci-003 | | | ChatGPT | | | BARD | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot | 0-shot | 1-shot | 3-shot |
| De. | bAbI-15 | ✓ | 63.00 | 56.00 | 43.00 | 22.60 | 39.40 | 55.70 | 99.00 | 84.00 | 62.00 |
| | EntailmentBank | ✓ | 8.00 | 6.00 | 7.00 | 7.06 | 5.88 | 3.24 | 26.00 | 25.00 | 28.00 |
| | RuleTaker | | 26.00 | 29.00 | 27.00 | 21.30 | 27.80 | 34.80 | 80.00 | 84.00 | 75.00 |
| | FOLIO | | 14.00 | 23.00 | 21.00 | 31.86 | 22.55 | 19.61 | 60.00 | 63.00 | 68.00 |
| | Leap-Of-Thought | | 71.00 | 55.00 | 54.00 | 32.74 | 5.04 | 4.73 | 2.00 | 2.00 | 0.00 |
| In. | bAbI-16 | ✓ | 60.00 | 77.00 | 86.00 | 93.60 | 29.80 | 41.20 | 96.00 | 98.00 | 99.00 |
| | CLUTRR | ✓ | 2.00 | 28.00 | 31.00 | 2.62 | 1.57 | 0.87 | 2.00 | 6.00 | 14.00 |
| Ab. | $\alpha$-NLI | | 2.00 | 2.00 | 1.00 | 1.00 | 0.20 | 0.10 | 8.00 | 16.00 | 0.00 |
| | $\alpha$-NLG | ✓ | 63.00 | 61.00 | 72.00 | 70.70 | 69.70 | 64.50 | 24.00 | 32.00 | 31.00 |
| | AbductiveRules | ✓ | 1.00 | 0.00 | 0.00 | 42.40 | 5.40 | 0.50 | 67.00 | 48.00 | 22.00 |
| | D*-Ab | ✓ | 85.00 | 27.00 | 17.00 | 55.30 | 27.10 | 16.70 | 18.00 | 16.00 | 2.00 |
| Mix | ReClor | | 1.00 | 1.00 | 1.00 | 2.00 | 1.20 | 1.40 | 11.00 | 16.00 | 24.00 |
| | LogiQA | | 0.00 | 5.00 | 0.00 | 1.54 | 0.77 | 1.08 | 32.00 | 35.00 | 43.00 |
| | LogiQA 2.0 | | 0.00 | 0.00 | 0.00 | 0.80 | 4.00 | 0.8 | 5.00 | 5.00 | 4.00 |
| | LogiQA2NLI | | 0.00 | 0.00 | 0.00 | 0.17 | 0.50 | 0.17 | 11.00 | 31.00 | 4.00 |

Further, LLMs are best at keeping rigorous reasoning in the abductive setting, while they are weak in the deductive and inductive settings. The finding is a little different from the analysis of simple accuracy conditions in the previous section. We argue that the setting of abductive reasoning requires the LLMs to achieve the reasoning reversely, which can activate LLMs to provide sufficient reasoning process. While in a deductive reasoning setting, the reasoning chain is sequential, which may cause LLMs to be in lazy mode and harm rigorous reasoning.

In Table 2, we also include the two conditions (1) when correct answer and correct explanations are satisfied, and (2) when correct answer and complete explanations are satisfied. Results vary a lot with different datasets and different LLMs. Overall, ChatGPT performs relatively well in keeping correct explanations while it may fail to maintain complete explanations in most cases. In comparison, text-davinci-003 exhibits stronger characteristics in maintaining the completeness of explanations, compared with the correctness of explanations. These findings of the respective reasoning preferences are expected to guide the future utilization of LLMs.

### 5.2 Are LLMs Self-aware Logical Reasoners?

From an alternative perspective, the redundancy of the generated content by LLMs has been a frequently discussed topic, as it is deemed an important metric for assessing their practicality. In this paper, we consider LLMs with less redundant content as more *self-aware*, as they can effectively express the necessary information without outputting all possible answers. Table 3 presents the evaluation results of LLMs' self-awareness. Similar to our previous approach, we compute the weighted results for each reasoning setting and derive the self-awareness scores shown in Fig. 5b. The darker color indicates a stronger self-awareness capability. Results indicate that text-davinci-003 exhibits notable advantages, particularly in the inductive, abductive, and mixed-form reasoning settings. Additionally, it ranks second in the deductive setting. Conversely, BARD performs poorly in deductive, abductive, and mixed-form reasoning settings.

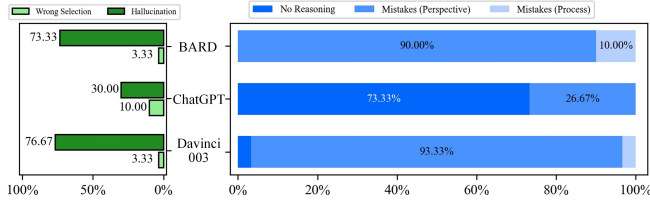In comparison to classification tasks, LLMs tend to generate redundant answers more frequently in generation tasks, such as $\alpha$-NLI vs. $\alpha$-NLG. This is because open-ended questions can prompt LLMs to generate content from various perspectives, which can lead to the inclusion of redundant information. Furthermore, the mixed-form reasoning setting observes significantly fewer instances of redundancy. The tasks in mixed-form reasoning are primarily based on question answering, which closely resembles real-life text, and LLMs tend to generate rational and specific content in such scenarios. However, in other settings, the input context is elaborately designed for logical reasoning and may provide sufficient background information. This can result in LLMs employing embodied commonsense knowledge to help reason and thus generate additional explanations.
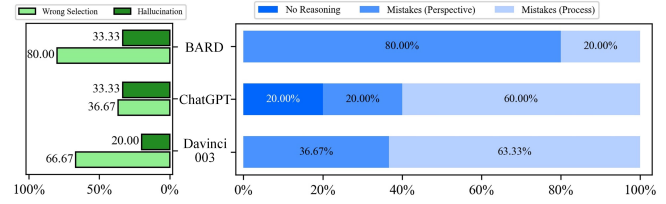
### 5.3 Do LLMs Have Obvious Logical Flaws?

Based on the preceding statements, we establish error types for problematic cases (with incorrect explanations) from two dimensions: *evidence selection process* and *reasoning process*. The former dimension encompasses two error types: (i) *wrong selection* and (ii) *hallucination*, which are independent of each other. The latter dimension comprises three error types: (i) *no reasoning*, (ii) *perspective mistake* and (iii) *process mistake*. Each problematic case can only be attributed to one of the three errors in the *reasoning process* dimension.

We visualize the attribution results for fourteen datasets in Fig. 6, including four deductive, two inductive, four abductive and four mixed-form ones. Overall speaking, the types of errors vary between datasets. For the *wrong selection*, 33.26% of the problematic cases fail to select the right answers for reasoning. Also, 27.46% suffers from the hallucination issue of LLMs. From the dimension of *reasoning process*, *no reasoning* error keeps a small portion in most of the cases, only covering 19.33% of the selected cases in total. Meanwhile, *perspective mistake* occupies 44.47% of the cases and *process mistake* covers 36.20%.
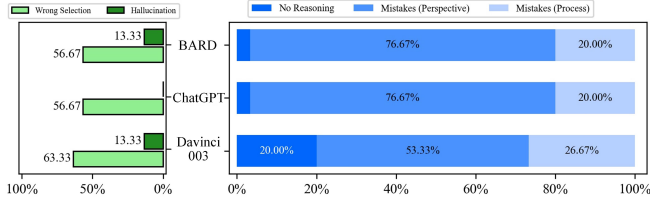
According to the statistics, the most obvious obstacle to logical reasoning tasks in LLMs is whether they can find the correct evidence and perspective. Due to the limited content of inputs, LLMs are also prone to generate hallucinatory facts to aid the reasoning process, which may affect the reliability of LLMs in real applications. In addition, LLMs abandon reasoning in a considerable number of cases. Such
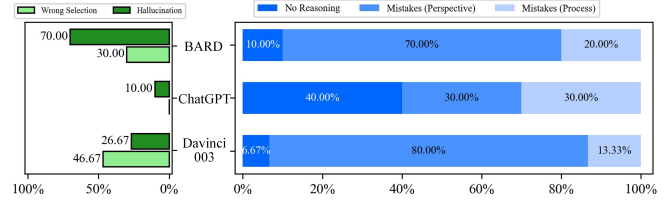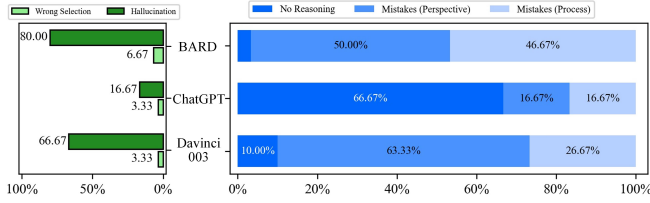
Fig. 6: Statistics of different error types from *evidence selection process* and *reasoning process* view.

|        | De. | In. | Ab. | Mix |
|--------|-----|-----|-----|-----|
| Davinci | 0.58 | 0.68 | 0.05 | 0.38 |
| ChatGPT | 0.21 | 0.00 | 0.52 | 0.36 |
| BARD | 0.91 | 0.88 | 0.59 | 0.37 |

(a) Visualization on the activity.

|        | De. | In. | Ab. | Mix |
|--------|-----|-----|-----|-----|
| Davinci | 0.26 | 0.06 | 0.23 | 0.26 |
| ChatGPT | 0.57 | 0.78 | 0.20 | 0.01 |
| BARD | 0.02 | 0.11 | 0.30 | 0.37 |

(b) Visualization on the orientation.

|        | De. | In. | Ab. | Mix |
|--------|-----|-----|-----|-----|
| Davinci | 0.32 | 0.40 | 0.47 | 0.56 |
| ChatGPT | 0.61 | 0.83 | 0.46 | 0.85 |
| BARD | 0.01 | 0.00 | 0.14 | 0.13 |

(c) Visualization on the no-hallucination.

Fig. 7: Heatmap results for the activity, orientation and no-hallucination of LLMs.



(a) EntailmentBank (De.).
(b) RuleTaker (De.).
(c) FOLIO (De.).
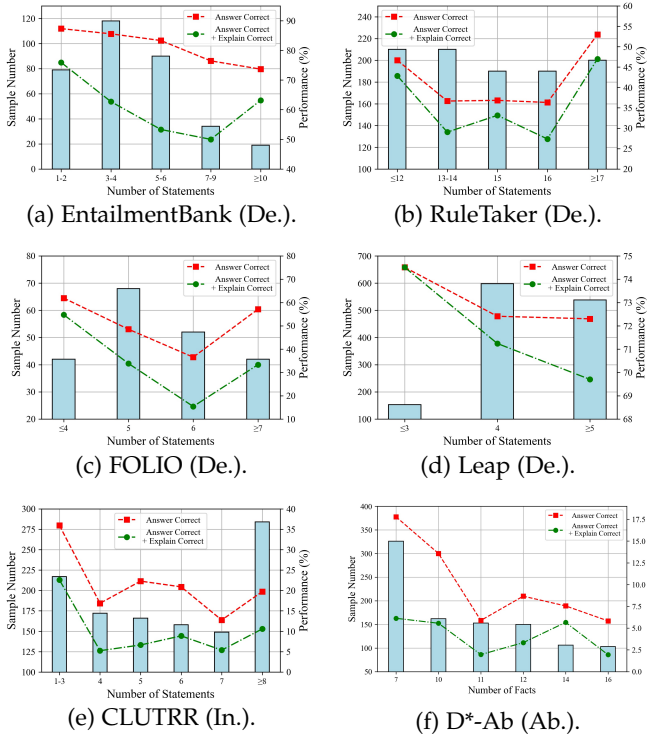(d) Leap (De.).
(e) CLUTRR (In.).
(f) D*-Ab (Ab.).

Fig. 8: The LLM performances with different numbers of statements.

phenomena of laziness are worth noting, especially when we depend on LLMs to help reason in the downstream tasks. In the following section, we will provide a detailed analysis of specific LLM and specific reasoning settings. Considering the above-mentioned obvious logical flaws, we will especially focus on the detailed analysis of the activity in reasoning, the orientation selection of LLMs and model hallucination.

### 5.4 Are LLMs Active Logical Reasoners?

In this paper, we consider LLMs with fewer *no reasoning* errors as more *Active* logical reasoners. Fig. 7a presents the weighted results for measuring active reasoning cases, where higher values indicate LLMs are more active in reasoning, while lower values represent lazier cases. Among the three LLMs, BARD is the most active logical reasoner, excelling in deductive, inductive and abductive settings. ChatGPT, on the other hand, is deemed the lazier reasoner

in deductive and inductive settings, while text-davinci-003 is lazier in abductive reasoning tasks.

Furthermore, we compare the performances of LLMs across different reasoning modes. In deductive reasoning tasks, LLMs exhibit more active reasoning, while in abductive settings, they tend to display lazier performances. Deductive tasks are in a forward reasoning mode, which is more natural for both generative LLMs and humans. This can inspire LLMs to generate effective reasoning chains. Conversely, abductive reasoning requires LLMs to provide explanations for the given inputs, which is in a backward reasoning mode. It is intuitive that LLMs may struggle to conduct reasoning in some cases.

### 5.5 Are LLMs Oriented Logical Reasoners?

At the outset of the reasoning process, it is crucial to identify the correct starting points and potential directions for reasoning. We consider LLMs with this capability as *Oriented* logical reasoners and present evaluation results based on the error type of *perspective mistake* in Fig. 7b. From the heatmap results, ChatGPT exhibits better oriented capability in deductive and inductive tasks compared to the other two LLMs. However, it frequently fails to identify the correct reasoning direction in abductive and mixed-form reasoning. Conversely, BARD performs well in identifying the right direction in abductive and mixed-form settings but struggles in deductive and inductive ones. Compared to the others, text-davinci-003 displays moderate performance in identifying reasoning perspectives.

In light of the preceding findings, ChatGPT is a lazier logical reasoner, but it excels at identifying the correct direction for reasoning. In other words, ChatGPT is prone to conduct confident reasoning. Conversely, text-davinci-003 and BARD are more active in logical reasoning, but they tend to start from the wrong direction, leading to reasoning mistakes.

### 5.6 Are LLMs Easy to Induce Hallucination in Logical Reasoning?

In the typical definition, hallucination refers to generated content that contradicts commonsense or current facts. To align with logical reasoning tasks, we expand the definition to include cases where facts are employed that contradict the context or are not verified by the context. Fig. 7c displays the weighted performances of cases with no hallucinations. The darker color indicates better performance in avoiding hallucinations. Based on the results, ChatGPT exhibits strong
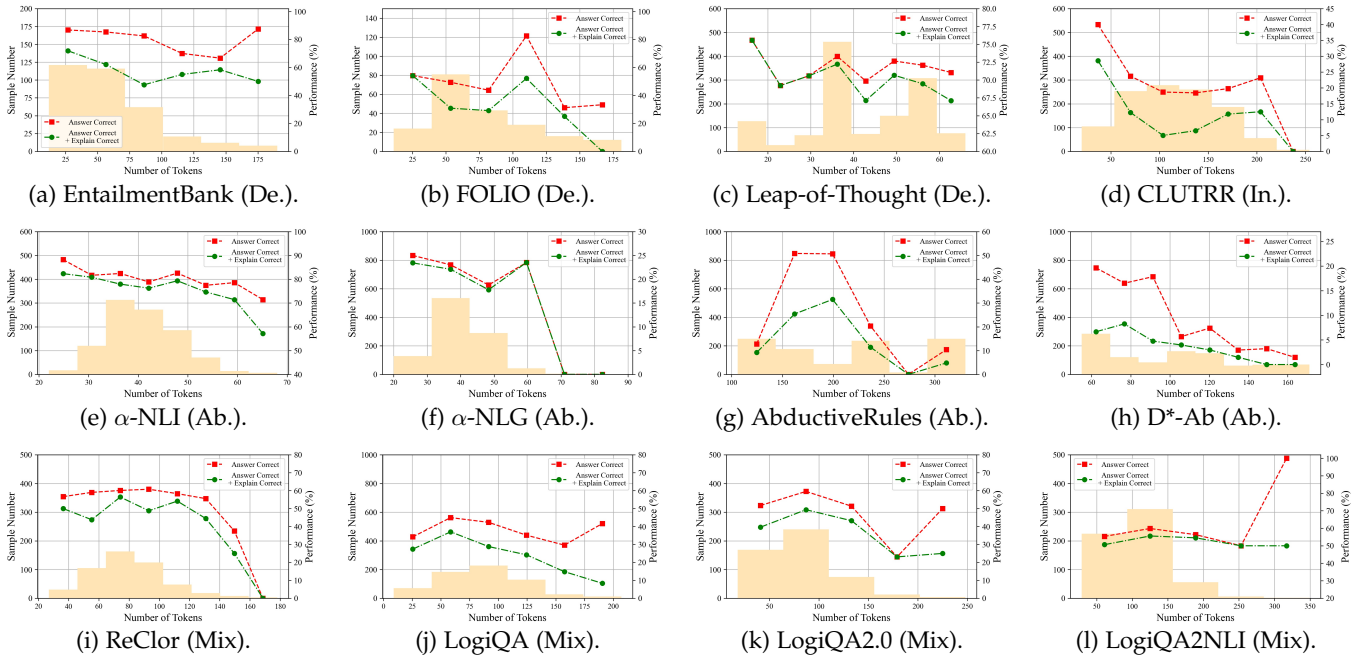
Fig. 9: The performances of ChatGPT with different tokens on various datasets.
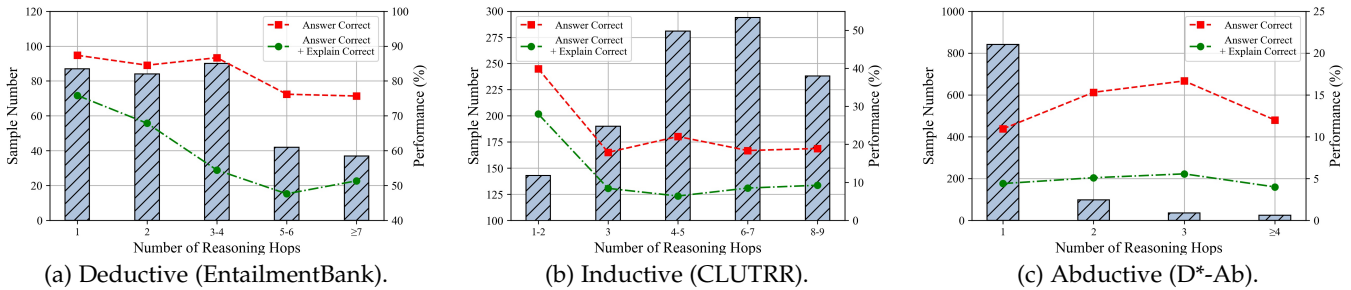


Fig. 10: The performances of ChatGPT under different number of hops. Comparison of Deductive, Inductive and Abductive reasoning settings.

and consistent competitiveness, ranking first in deductive, inductive, and mixed-form settings. It also ranks second in abductive reasoning tasks, albeit with slight disadvantages. Conversely, BARD displays poor performance in avoiding hallucinations and maintaining clarity during reasoning. Across all four reasoning settings, BARD ranks last with significant gaps.

On average, LLMs induce hallucination in 27.46% of the failure cases. Among deductive, inductive and abductive settings, model hallucinations are more common under deductive reasoning tasks. While LLMs may have clearer minds in the inductive setting.

### 5.7 How Does the Number of Statements Affect the LLMs' Performances?

In the following, we explore some of the key factors to affect the reasoning performances of LLMs. Since the length of the input context can be different for the datasets, we report the model performances with the number of statements in Fig. 8. We take ChatGPT for analysis and choose six datasets with specific counts of statements for illustration, covering

the three reasoning manners. The first four subfigures are related to the deductive reasoning manner, i.e., Entailment-Bank, RuleTaker, FOLIO, and Leap-of-Thought. Fig. 8e is CLUTRR in the inductive setting and Fig. 8f is D-Ab in the abductive setting. The horizontal axis denotes the number of statements. The left vertical axis denotes the number of samples for different numbers of statements. And the right vertical axis represents the performances with different numbers of statements.

From the overall results, LLMs can keep the correctness of both the answer and explanations with fewer input statements. With the statement number increasing, the performances drop a lot and LLMs struggle to give the correct explanations. Interestingly, five (out of six) datasets witness performance gains when the number of statements reaches certain values. For example, in the RuleTaker dataset, the best performances are achieved when the number of statements is larger than 17. And when the number is between 13 to 16, ChatGPT is capable of keeping stable performances. We argue that the larger number of statements can provide richer information and sometimes can help control the rea-

TABLE 4: Statistics and evaluation results on NeuLR. *Num.* represents the number of samples in the dataset. *#Hop* represents the hop number of samples in the dataset. *COT* represents the chain-of-thought strategy under the 1-shot setting.

| Dataset | Num. | #Hop | text-davinci-003 | | | ChatGPT | | | BARD | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 0-shot | 1-shot | COT | 0-shot | 1-shot | COT | 0-shot | 1-shot | COT |
| **NeuLR** | 3,000 | 1~5 | 50.93 | 59.17 | 67.90 | 37.27 | 48.13 | 48.00 | 63.67 | 65.07 | 66.00 |
| - Deductive | 1,000 | 2 | 59.00 | 69.40 | 86.10 | 85.20 | 69.10 | 68.30 | 87.40 | 93.10 | 91.90 |
| - Inductive | 1,000 | 3 | 86.90 | 89.60 | 95.60 | 15.10 | 68.60 | 69.60 | 96.00 | 92.60 | 96.30 |
| - Abductive | 1,000 | 1~5 | 6.90 | 18.50 | 22.00 | 11.50 | 6.70 | 6.10 | 7.60 | 9.50 | 9.80 |

soning direction of LLMs.

Furthermore, we investigate the impact of the number of tokens in the context. Fig. 9 illustrates the ChatGPT performances on twelve datasets. Results vary a lot across different datasets and different reasoning settings. In general, as the number of tokens increases, the performances of ChatGPT tend to decline. Detailedly, ChatGPT performs stably in deductive settings. Especially for EntailmentBank and Leap-of-Thought datasets, ChatGPT maintains relatively consistent accuracy with token numbers increasing. Considering that deductive setting is a more common form of reasoning in reality, LLMs are well-trained on it to tackle the various lengths of inputs.

## 5.8 How Does the Number of Reasoning Hops Affect the LLMs' Performances?

Also, it is interesting to explore the influences of reasoning hops for LLMs. Among the selected dataset, three of them offer the number of hops for each sample, which are EntailmentBank in deductive reasoning, CLUTRR in inductive reasoning and D*-Ab in abductive reasoning. Fig. 10 presents the performance of ChatGPT with different hop numbers (Results of other LLMs are listed in the Appendix). Alongside the simple accuracy results, we also report the rigorous reasoning cases where both the answer and explanations are correct.

In the deductive setting, with the number of hops increasing, the performances witness obvious drops, particularly influencing the rigor of the LLM reasoning. It illustrates that reasoning hops have great effects on deductive reasoning. In inductive reasoning, when the hop number is greater than two, the performance of ChatGPT decreases sharply. When the number ranges from three to nine, the performance of ChatGPT keeps stable at a relatively low level. Combined with the weak performance of ChatGPT in inductive reasoning tasks, it demonstrates that ChatGPT can only work on simple induction, and it obviously struggles in cases when more hops are needed. In the abductive reasoning setting, the majority of the test samples only need one-hop reasoning. When the hop number increases, the performances of ChatGPT witness slight improvements. It shows that ChatGPT may have the potential capability of multi-hop reasoning in the abductive setting.

## 6 NEUTRAL-CONTENT LOGICAL REASONING

Considering the current benchmarks may not provide neutral content for fair evaluation, we propose the new dataset NeuLR to benchmark the neutral-content logical reasoning tasks. In column 1~3 of TABLE 4, we provide the statistics of NeuLR. It contains 3k samples in total, with 1k for deductive reasoning, 1k for inductive reasoning and 1k for abductive reasoning. Limited by space, we provide the details of the construction of NeuLR in the Appendix.

To evaluate the performances of LLMs on NeuLR, we conduct the experiments shown in TABLE 4. Especially, we provide three different test settings, i.e., zero-shot, one-shot and chain-of-thought [48] settings. Details of the prompt forms are included in the Appendix.

Firstly, from the results, few-shot prompting and chain-of-thought prompting can both boost the performances of LLMs in most cases. Overall, chain-of-thought helps most to the model accuracy. Especially for text-davinci-003, it witnesses consistent gains with the aid of few-shot prompting and chain-of-thought prompting strategies.

Secondly, among the zero-shot results of three LLMs, BARD achieves the best performances on NeuLR while ChatGPT ranks last. The differences of zero-shot settings are significant. However, with the help of few-shot and chain-of-thought prompting strategies, text-davinci-003 large narrows the gaps with BARD, and it surpasses BARD with chain-of-thought strategy. Overall, the performances of LLMs on NeuLR still have great room for improvement.

Thirdly, from the perspective of different reasoning settings, there exist huge differences in results compared with the findings in previous sections. Generally speaking, the LLMs' performances on inductive reasoning are better than deductive reasoning and abductive reasoning. Especially, LLMs present obvious weakness in the abductive setting. While from fifteen classical datasets, the performances among the reasoning settings are sorted as: *deductive > abductive > inductive*. Such findings can also motivate future studies. For example, content neutrality may help design related strategies to improve the performance of inductive reasoning.

## 7 CONCLUSION

In this paper, in-depth evaluations are conducted on logical reasoning tasks, discussing whether LLMs are really good logical reasoners. First, the logical reasoning evaluations are organized from deductive, inductive, abductive and mixed-form views. We select fifteen logical reasoning datasets to evaluate on three representative LLMs (i.e., text-davinci-003, ChatGPT and BARD) under both zero-shot and few-shot settings. Second, this paper provides fine-level evaluations on four metrics, covering both objective and subjective views. For problematic cases, extensive error attributions are conducted from two dimensions, forming five error types. It uncovers the logical flaws of LLMs and we provide deep analysis on the results. Third, to achieve a fair and pure benchmark for logical reasoning capability, we propose a
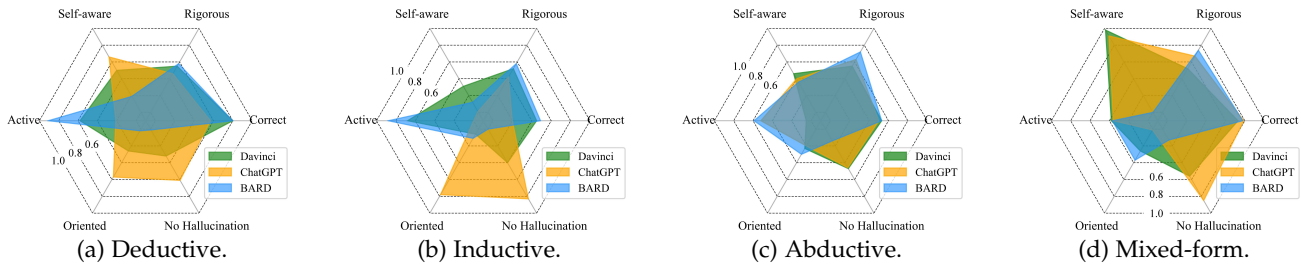
Fig. 11: Visualization of LLM capability under four reasoning settings.

dataset with neutral content, covering deductive, inductive and abductive settings.

Based on the evaluation results above, we abstract six dimensions to measure the logical reasoning capability of LLMs: (1) *Correct*, (2) *Rigorous*, (3) *Self-aware*, (4) *Active*, (5) *Oriented* and (6) *No hallucination*. All these properties can be calculated with the evaluation methods proposed in this paper. Therefore, we propose an evaluation scheme for the logical reasoning capability of LLMs. Considering the different performances of LLMs on deductive, inductive, abductive and mixed-form settings, we respectively visualize each ability map in Fig. 11.

According to the results, text-davinci-003 can maintain balanced performances in deductive and mixed-form settings. But it usually fails to keep oriented for reasoning in the inductive setting, and it also shows laziness in the abductive reasoning tasks. Since it is the earliest released LLM of the three, it is understandable that text-davinci-003 has some limitations in logical reasoning tasks, especially in the more complex settings (e.g., inductive and abductive).

From the perspective of common evaluations, ChatGPT is the weakest LLM of the three, since it performs badly in showing correct and rigorous reasoning under deductive, inductive and abductive settings. Also, it seems to be the laziest reasoner in deductive and inductive settings. However, it surprises us that it shows unique advantages in maintaining oriented reasoning and avoiding hallucination, especially in deductive and inductive settings. In addition, it shows its comprehensive capability in the mixed-form setting. We argue that ChatGPT is specially designed for chatting, thus it does pretty well in keeping rational but is not good at solving complex reasoning problems.

BARD is the most active reasoner and it keeps great competitiveness as a correct and rigorous reasoner. However, it also shows obvious flaws compared with other LLMs. BARD tends to generate redundant content, easily fails to find the correct reasoning directions and it usually fails to avoid hallucinations. In short, BARD shows great advantages in current benchmarks with objective metrics, due to the larger model size and massive training data. But it still has much room for improvement in some implicit aspects, i.e., self-awareness, orientation and non-hallucination.

Overall, it can be observed that all LLMs exhibit specific limitations in logical reasoning, with relative strength in deductive reasoning but evident struggles in inductive settings. Moreover, current evaluation benchmarks, which primarily depend on objective metrics, are not sufficient to comprehensively evaluate LLMs.

## 8 FUTURE DIRECTIONS

Based on the evaluation results, this paper concludes six future directions for logical reasoning tasks.

**Strengthen the reasoning ability of inductive reasoning.** Inductive reasoning draws broad conclusions from specific observations, requiring a more abstract and comprehensive understanding of real-world knowledge compared to deductive or abductive reasoning. However, LLMs have shown poor performance in this area, as demonstrated in Section 4.1. Therefore, it is crucial to develop pre-training or fine-tuning strategies to enhance their inductive reasoning abilities. One such strategy could be constructing more inductive instructions to guide LLMs.

**Enhance the LLM's perception of its capability boundaries.** LLMs are capable of generating answers and explanations for reasoning questions regardless of difficulty and rationality. To realize it, LLMs would list some irrelevant facts of the given context or even hallucinations, as Sections 5.1 and 5.4 demonstrate. It will lead to LLMs solemnly talking nonsense and resulting in illogical, uninformative, or meaningless answers. A good logical reasoner should be aware of its boundaries and acknowledge when it is unable to answer a question. To enhance LLMs' self-awareness of their capability boundaries, future research could focus on cognitive science and neuroscience studies of human self-awareness.

**Strengthen the rigorous reasoning to apply to real-world scenarios.** Table 2 illustrates that current LLMs are not sufficiently rigorous for deductive, inductive, abductive, and mixed reasoning. As a result, there is still a significant gap between their capabilities and their potential applications in real-world scenarios, particularly those that require detailed intermediate explanations. For instance, using LLMs to solve mathematical problems and provide precise intelligent Q&A services in the education field remains a significant challenge [49].

**Minimize the occurrence of hallucinations.** Similar to the behaviors in other problem-solving contexts [50], LLMs may generate false or irrelevant hallucinations during logical reasoning tasks. It suggests that LLMs may not fully comprehend the question and can not solve it correctly. To address this issue, future research should develop more comprehensive evaluation metrics for hallucinations and explore specific strategies to minimize their occurrence.

**Improve the multi-hop reasoning capability, especially in inductive and abductive settings.** Combined with the results from Fig. 10 and Fig. 1 in Appendix, the multi-hop reasoning capability of LLMs still have much room for improvement. Especially in the inductive and abductive

settings, LLMs perform quite struggling. Since the multi-hop reasoning evaluates the high-level capabilities of LLMs, it is necessary to extend LLM capability to such complex settings. In fact, humans are better at decomposing complex questions. It can be an interesting topic for LLMs to capture the ability to divide and conquer questions, thus benefiting multi-hop reasoning.

**Increase explainability.** Finally, the explainability of LLMs will be essential for building trust, detecting and mitigating biases, improving performance, promoting user understanding, and complying with regulations. A commonsense-based neuro-symbolic AI framework, such as the one proposed by [51] for sentiment analysis, can help increase the explainability of the reasoning processes required for decision-making, which is crucial for sensitive applications involving ethics, privacy and health.

## ACKNOWLEDGMENTS

## REFERENCES

[1] G. Antoniou and A. Bikakis, "Dr-prolog: A system for defeasible reasoning with rules and ontologies on the semantic web," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, vol. 19, no. 2, pp. 233–245, 2007.

[2] T. Lukasiewicz, "A novel combination of answer set programming with description logics for the semantic web," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, vol. 22, no. 11, pp. 1577–1592, 2010.

[3] Q. Lin, J. Liu, L. Zhang, Y. Pan, X. Hu, F. Xu, and H. Zeng, "Contrastive graph representations for logical formulas embedding," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, vol. 35, no. 4, pp. 3563–3574, 2023.

[4] L. Wu, Y. Zhou, and D. Zhou, "Towards high-order complementary recommendation via logical reasoning network," in *IEEE International Conference on Data Mining (ICDM)*, 2022, pp. 1227–1232.

[5] Q. Lin, J. Liu, F. Xu, Y. Pan, Y. Zhu, L. Zhang, and T. Zhao, "Incorporating context graph with logical reasoning for inductive relation prediction," in *The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2022, pp. 893–903.

[6] J. Yu, Q. Su, X. Quan, and J. Yin, "Multi-hop reasoning question generation and its application," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, vol. 35, no. 1, pp. 725–740, 2023.

[7] H. Bronkhorst, G. Roorda, C. Suhre, and M. Goedhart, "Logical reasoning in formal and everyday reasoning tasks," *International Journal of Science and Mathematics Education*, vol. 18, pp. 1673–1694, 2020.

[8] Z. Yang, Z. Du, R. Mao, J. Ni, and E. Cambria, "Logical reasoning over natural language as knowledge representation: A survey," *CoRR*, vol. abs/2303.12023, 2023.

[9] B. Peng, C. Li, P. He, M. Galley, and J. Gao, "Instruction tuning with GPT-4," *CoRR*, vol. abs/2304.03277, 2023. [Online]. Available: https://doi.org/10.48550/arXiv.2304.03277

[10] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J. Nie, and J. Wen, "A survey of large language models," *CoRR*, vol. abs/2303.18223, 2023.

[11] M. M. Amin, E. Cambria, and B. W. Schuller, "Will affective computing emerge from foundation models and general artificial intelligence? A first evaluation of chatgpt," *IEEE Intelligent Systems*, vol. 38, no. 2, pp. 15–23, 2023.

[12] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung, Q. V. Do, Y. Xu, and P. Fung, "A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity," *CoRR*, vol. abs/2302.04023, 2023.

[13] N. Bian, X. Han, L. Sun, H. Lin, Y. Lu, and B. He, "Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models," *CoRR*, vol. abs/2303.16421, 2023.

[14] S. Imani, L. Du, and H. Shrivastava, "Mathprompter: Mathematical reasoning using large language models," *CoRR*, vol. abs/2303.05398, 2023.

[15] P. A. Flach and A. C. Kakas, "Abductive and inductive reasoning: background and issues," *Abduction and induction: Essays on their relation and integration*, pp. 1–27, 2000.

[16] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe, "Training language models to follow instructions with human feedback," in *NeurIPS*, 2022.

[17] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.

[18] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen *et al.*, "Palm 2 technical report," *arXiv preprint arXiv:2305.10403*, 2023.

[19] C. Qin, A. Zhang, Z. Zhang, J. Chen, M. Yasunaga, and D. Yang, "Is chatgpt a general-purpose natural language processing task solver?" *CoRR*, vol. abs/2302.06476, 2023.

[20] S. Tu, C. Li, J. Yu, X. Wang, L. Hou, and J. Li, "Chatlog: Recording and analyzing chatgpt across time," *CoRR*, vol. abs/2304.14106, 2023.

[21] Y. Huang, Y. Bai, Z. Zhu, J. Zhang, J. Zhang, T. Su, J. Liu, C. Lv, Y. Zhang, J. Lei, F. Qi, Y. Fu, M. Sun, and J. He, "C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models," *CoRR*, vol. abs/2305.08322, 2023.

[22] T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. R. McKeown, and T. B. Hashimoto, "Benchmarking large language models for news summarization," *CoRR*, vol. abs/2301.13848, 2023.

[23] Y. Pan, J. Liu, L. Zhang, T. Zhao, Q. Lin, X. Hu, and Q. Wang, "Inductive relation prediction with logical reasoning using contrastive representations," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2022, pp. 4261–4274.

[24] Z. Huang, M. Chiang, and W. Lee, "Line: Logical query reasoning over hierarchical knowledge graphs," in *28th SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 615–625.

[25] K. Cheng, J. Liu, W. Wang, and Y. Sun, "Rlogic: Recursive logical rule learning from knowledge graphs," in *28th SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 179–189.

[26] P. N. Johnson-Laird, "Deductive reasoning," *Annual review of psychology*, vol. 50, no. 1, pp. 109–135, 1999.

[27] V. Goel, "Anatomy of deductive reasoning," *Trends in cognitive sciences*, vol. 11, no. 10, pp. 435–441, 2007.

[28] E. Heit and C. M. Rotello, "Relations between inductive reasoning and deductive reasoning." *Journal of Experimental Psychology: Learning, Memory, and Cognition*, vol. 36, no. 3, p. 805, 2010.

[29] F. Yu, H. Zhang, and B. Wang, "Nature language reasoning, a survey," *arXiv preprint arXiv:2303.14725*, 2023.

[30] J. R. Josephson and S. G. Josephson, *Abductive inference: Computation, philosophy, technology*. Cambridge University Press, 1996.

[31] D. Walton, "Abductive, presumptive and plausible arguments," *Informal Logic*, vol. 21, no. 2, 2001.

[32] B. Dalvi, P. Jansen, O. Tafjord, Z. Xie, H. Smith, L. Pipatanangkura, and P. Clark, "Explaining answers with entailment trees," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021, pp. 7358–7370.

[33] S. Han, H. Schoelkopf, Y. Zhao, Z. Qi, M. Riddell, L. Benson, L. Sun, E. Zubova, Y. Qiao, M. Burtell, D. Peng, J. Fan, Y. Liu, B. Wong, M. Sailor, A. Ni, L. Nan, J. Kasai, T. Yu, R. Zhang, S. R. Joty, A. R. Fabbri, W. Kryscinski, X. V. Lin, C. Xiong, and D. Radev, "FOLIO: natural language reasoning with first-order logic," *CoRR*, vol. abs/2209.00840, 2022.

[34] A. Talmor, O. Tafjord, P. Clark, Y. Goldberg, and J. Berant, "Leap-of-thought: Teaching pre-trained models to systematically reason over implicit knowledge," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

[35] K. Sinha, S. Sodhani, J. Dong, J. Pineau, and W. L. Hamilton, "CLUTRR: A diagnostic benchmark for inductive reasoning from text," in *Proceedings of the 2019 Conference on Empirical Methods in*

*Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 4505–4514.

[36] W. Yu, Z. Jiang, Y. Dong, and J. Feng, "Reclor: A reading comprehension dataset requiring logical reasoning," in *8th International Conference on Learning Representations (ICLR)*, 2020.

[37] J. Liu, L. Cui, H. Liu, D. Huang, Y. Wang, and Y. Zhang, "Logiqa: A challenge dataset for machine reading comprehension with logical reasoning," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI)*, 2020, pp. 3622–3628.

[38] H. Liu, R. Ning, Z. Teng, J. Liu, Q. Zhou, and Y. Zhang, "Evaluating the logical reasoning ability of chatgpt and GPT-4," *CoRR*, vol. abs/2304.03439, 2023.

[39] J. Weston, A. Bordes, S. Chopra, and T. Mikolov, "Towards ai-complete question answering: A set of prerequisite toy tasks," in *4th International Conference on Learning Representations (ICLR)*, 2016.

[40] P. Clark, O. Tafjord, and K. Richardson, "Transformers as soft reasoners over language," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI)*, 2020, pp. 3882–3890.

[41] C. Bhagavatula, R. L. Bras, C. Malaviya, K. Sakaguchi, A. Holtzman, H. Rashkin, D. Downey, W. Yih, and Y. Choi, "Abductive commonsense reasoning," in *8th International Conference on Learning Representations (ICLR)*, 2020.

[42] N. Young, Q. Bao, J. Bensemann, and M. Witbrock, "Abduction-rules: Training transformers to explain unexpected inputs," in *Findings of the Association for Computational Linguistics*, 2022, pp. 218–227.

[43] O. Tafjord, B. Dalvi, and P. Clark, "Proofwriter: Generating implications, proofs, and abductive statements over natural language," in *Findings of the Association for Computational Linguistics*, 2021, pp. 3621–3634.

[44] P. Minervini, S. Riedel, P. Stenetorp, E. Grefenstette, and T. Rocktäschel, "Learning reasoning strategies in end-to-end differentiable proving," in *Proceedings of the 37th International Conference on Machine Learning (ICML)*, ser. Proceedings of Machine Learning Research, vol. 119, 2020, pp. 6938–6949.

[45] F. Jiao, Y. Guo, X. Song, and L. Nie, "Merit: Meta-path guided contrastive learning for logical reasoning," in *Findings of the Association for Computational Linguistics (Findings of ACL)*, 2022, pp. 3496–3509.

[46] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, "A survey for in-context learning," *arXiv preprint arXiv:2301.00234*, 2022.

[47] Z. Wu, Y. Wang, J. Ye, and L. Kong, "Self-adaptive in-context learning," *arXiv preprint arXiv:2212.10375*, 2022.

[48] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *NeurIPS*, 2022.

[49] I. Drori, S. Zhang, R. Shuttleworth, L. Tang, A. Lu, E. Ke, K. Liu, L. Chen, S. Tran, N. Cheng *et al.*, "A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level," *Proceedings of the National Academy of Sciences (PNAS)*, vol. 119, no. 32, p. e2123433119, 2022.

[50] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. Bang, A. Madotto, and P. Fung, "Survey of hallucination in natural language generation," *ACM Computing Surveys*, vol. 55, no. 12, pp. 248:1–248:38, 2023.

[51] E. Cambria, Q. Liu, S. Decherchi, F. Xing, and K. Kwok, "SenticNet 7: A commonsense-based neurosymbolic AI framework for explainable sentiment analysis," in *LREC*, 2022, pp. 3829–3839.