

# A Survey of Inductive Reasoning for Large Language Models

Anonymous ACL submission

## Abstract

Reasoning is an important task for large language models (LLMs). Among all the reasoning paradigms, inductive reasoning is one of the basic types, which is characterized by its particular-to-general thinking process and the non-uniqueness of its answers. The inductive mode is crucial for knowledge generalization and aligns better with human cognition, so it is a fundamental mode of learning, hence attracting increasing interest. Despite the importance of inductive reasoning, there is no systematic summary of it. Therefore, this paper presents the first comprehensive survey of inductive reasoning for LLMs. First, methods for improving inductive reasoning are categorized into three main areas: post-training enhancement, test-time exploration, and data augmentation. Then, current benchmarks of inductive reasoning are summarized, and a unified sandbox-based evaluation approach with the observation coverage metric is derived. Finally, we offer some analyses regarding the source of inductive ability and how simple model architectures and data help with inductive tasks, providing a solid foundation for future research.

## 1 Introduction

In recent years, the rapid development of large language models (LLMs) (Zhao et al., 2023) leads to significant progress in many natural language processing (NLP) downstream tasks. Among these, reasoning (Huang and Chang, 2022; Plaat et al., 2024; Zhang et al., 2025a) is one of the comprehensive and challenging tasks for LLMs, and therefore receives considerable attention.

Within all the reasoning paradigms, inductive reasoning (Lu, 2024) is one of the fundamental types. It involves drawing general conclusions from specific observations (Han et al., 2024). We give two examples illustrating number sequence calculation and list transformation in Figure 1. The main characteristics of inductive reasoning are

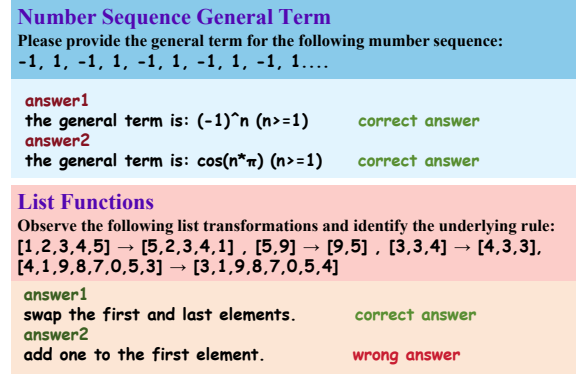


Figure 1: Two examples of inductive reasoning. They generalize from specific observations or cases to derive general conclusions or rules. There may be more than one such conclusion that meets all the observations.

its particular-to-general thinking process and the non-uniqueness of its answers. Considering how humans perceive the world, they typically make judgments by drawing analogies from past experiences to current situations, rather than always going through a strictly logical process as in deductive reasoning (Ingold, 2021). We can assume that the inductive mode is key to knowledge generalization and better aligns with human cognition. It is a fundamental mode of learning and thus attracts increasing interest.

Despite the importance of inductive reasoning, previous works mostly focus on deductive reasoning (Li et al., 2024b), represented by mathematical proof (Ahn et al., 2024; Chen et al., 2024b) and program verification (Liu et al., 2023; Jiang et al., 2024b), which is a logical reasoning that derives necessary conclusions from general rules or premises. For more conceptual distinctions, please refer to Appendix A.2. It has already been extensively studied in recent years (Lu et al., 2024; Wang et al., 2024a). Moreover, there is no systematic summary of inductive reasoning for LLMs.

Therefore, this paper presents the first compre-

067        hensive survey of inductive reasoning for LLMs.        114  
068        We introduce the background, including relevant        115  
069        concepts, applications in NLP and real-world scen-        116  
070        arios, as well as the significance (Section 2) at        117  
071        the beginning. The main body consists of three        118  
072        parts. First, we review and prospect the methods        119  
073        for enhancing the inductive reasoning capabilities        120  
074        of LLMs (Section 3), which are categorized into        121  
075        three main areas: post-training enhancement, test-        122  
076        time exploration, and data augmentation. We then        123  
077        summarize the current benchmarks of inductive        124  
078        reasoning and derive a unified sandbox-based eval-        125  
079        uation approach with the observation coverage met-        126  
080        ric (Section 4). Finally, we offer some theoretical        127  
081        analyses of inductive reasoning (Section 5), regard-        128  
082        ing the sources of inductive ability and practical        129  
083        experiences for enhancing it. The taxonomy of        130  
084        this survey is in Figure 2. In summary, the main        131  
085        contributions of this paper are threefold:        132

- 086        • **First survey.** To our knowledge, we are the        133  
087        first to present a comprehensive survey of induc-        134  
088        tive reasoning for LLMs, thoroughly analyz-        135  
089        ing the current techniques and applications.        136
- 090        • **New taxonomy.** We categorize methods for        137  
091        improving inductive reasoning into: post-        138  
092        training enhancement, test-time exploration,        139  
093        and data augmentation. We also summarize        140  
094        the current benchmarks and derive a unified        141  
095        sandbox-based evaluation approach.        142
- 096        • **Bright prospect.** We offer some analyses        143  
097        regarding the source of inductive ability and        144  
098        how simple model architectures and data help        145  
099        with inductive tasks, providing a solid founda-        146  
100        tion for future research (Appendix A.8).        147

## 101        2 Background        148

102        In this section, we will introduce the relevant con-        149  
103        cepts of inductive reasoning, some of its applica-        150  
104        tion scenarios, and the significance of studying it.        151

### 105        2.1 Concepts        152

#### 106        2.1.1 Large Language Models        153

107        Since the transformer architecture (Vaswani et al.,        154  
108        2017) has become mainstream for language mod-        155  
109        els, the field of NLP has experienced rapid devel-        156  
110        opment (Wei et al., 2021). During 2017 to 2022,        157  
111        pretrained language models (PLMs) (Kalyan et al.,        158  
112        2021), which undergo two stages—pretraining and        159  
113        finetuning—such as the BERT (Devlin et al., 2019)

and T5 (Raffel et al., 2023) series, once dominated        114  
the entire field. From 2022, with the advent of        115  
ChatGPT-3.5 (Ye et al., 2023), the era of LLMs of-        116  
ficially begins. LLMs, with their massive number        117  
of parameters and unique training methods (Zhao        118  
et al., 2025), significantly improve the performance        119  
of NLP tasks (She et al., 2023; Xu et al., 2023;        120  
Plaat et al., 2024) and have a profound impact on        121  
various aspects of daily life (Gan et al., 2023; Zhou        122  
et al., 2023; Maatouk et al., 2023). Some well-        123  
known LLMs include the GPT series<sup>1</sup>, the Gemini        124  
series<sup>2</sup>, the Claude series<sup>3</sup>, and so on.        125

#### 126        2.1.2 Inductive Reasoning        127

Inductive reasoning represents *making an induc-*        128  
*tion from specific instances or observations to de-*        129  
*rive general rules and conclusions* (Arthur, 1994;        130  
Heit, 2000). From another perspective, it denotes        131  
one reasoning approach where the conclusion is not        132  
guaranteed with certainty, but instead supported        133  
only to a certain degree of probability (Copi et al.,        134  
2004). In other words, inductive reasoning may        135  
have more than one valid hypothesis that can ac-        136  
count for all the instances or observations, making        137  
its answer open (Thomas, 2003). To sum up ab-        138  
stractly, inductive reasoning is a thought process        139  
that proceeds from the particular to the general.        140

### 141        2.2 Applications of Inductive Reasoning        142

The core idea of inductive reasoning is inductive        143  
bias. It is a set of assumptions or prior conditions        144  
that a model or an individual relies on when en-        145  
countering unseen items (Caruana, 1993; Baxter,        146  
2000). There is no ‘universal’ bias in deep learning.        147  
Choosing an appropriate inductive bias for a spe-        148  
cific task is key to achieving success (Provost and        149  
Buchanan, 1995). We will discuss its applications        150  
in NLP tasks and real-world scenarios.        151

#### 152        2.2.1 NLP Downstream Tasks        153

Inductive reasoning is widely applied to improve        154  
the performance of NLP downstream tasks (Ap-        155  
pendix A.3). Some common practices include        156  
training models to learn inductive bias (Yang et al.,        157  
2025), constructing chains of thought (CoT) (Chen        158  
et al., 2025b) or summarizing rules to enhance in-        159  
terpretability (Xu and Yang, 2025), and leveraging        160  
intra-parameters implicit knowledge (Cheng et al.,        161  
2024b) for induction, and others. Such approaches        162

<sup>1</sup><https://chatgpt.com/overview>

<sup>2</sup><https://cloud.google.com/vertex-ai/generative-ai/docs>

<sup>3</sup><https://claude.ai>

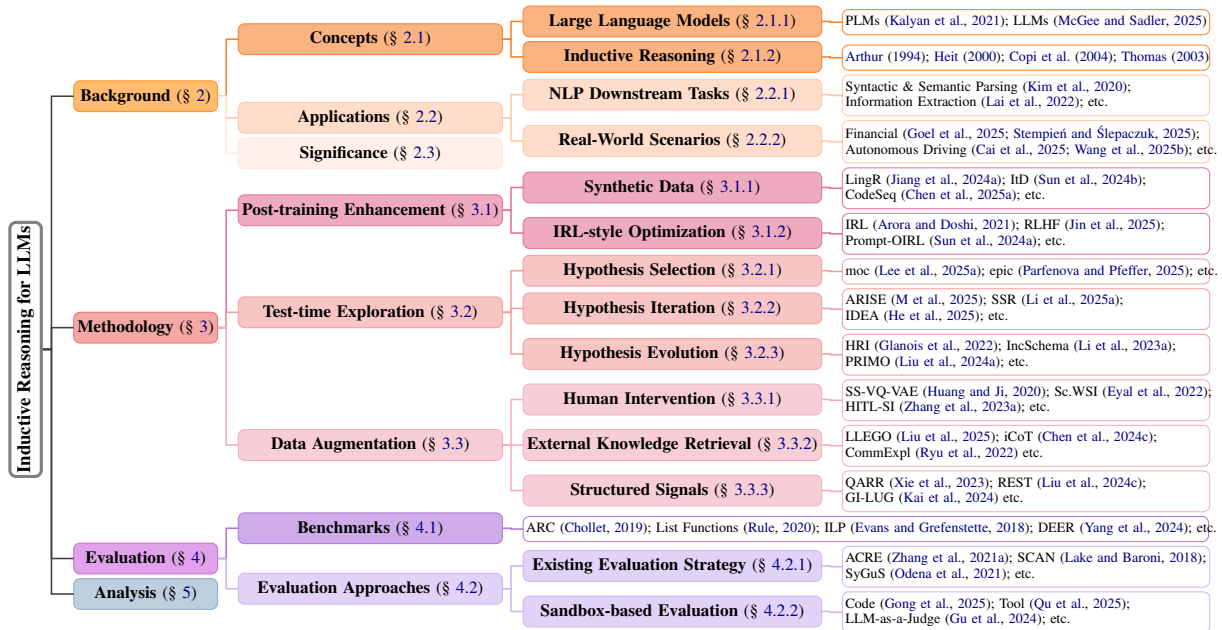


Figure 2: Taxonomy of the survey about the inductive reasoning for LLMs.

benefit a wide range of downstream NLP tasks: syntactic and semantic parsing (Kim et al., 2020; Yamada et al., 2021; Lindemann et al., 2024; Tsujimoto et al., 2025), information extraction (Lai et al., 2022; Liu et al., 2024c; Silva et al., 2025; Xu et al., 2024), dialogue systems (Feng et al., 2022; Xie et al., 2024; Ou et al., 2024), question answering (Gu and Su, 2022; Kim et al., 2023; Chen et al., 2024e), and multimodal tasks (Amosy et al., 2024; Zhou et al., 2025; Naik et al., 2025).

## 2.2.2 Real-World Scenarios

Inductive reasoning has a broad impact on real-world scenarios. We list three of them. (1) Financial forecasting: Inductive models are essential for predicting future financial outcomes by learning complex, non-linear patterns from vast amounts of historical time-series data (Faheem, 2021; Goel et al., 2025; Stempień and Ślepaczuk, 2025). (2) Autonomous driving: Inductive reasoning enables autonomous driving systems to generalize from historical data and past experiences to recognize rules in traffic conditions to make decisions. (Cai et al., 2025; Wang et al., 2025b). (3) Conversational healthcare and diagnostic dialogue: Inductive reasoning empowers artificial intelligence systems to mimic a clinician’s process of taking patient history and formulating a diagnosis by generalizing from symptom patterns (Tu et al., 2024; Dhudum et al., 2024; Zhang et al., 2025b). Broader real-world applications are in Appendix A.8.

## 2.3 Significance of Inductive Reasoning

Inductive reasoning has broad applications in both AI and real-world scenarios. It is the most universal and essential method in knowledge discovery and generalization (Carter and Hamilton, 1998; Bai et al., 2024; Sun et al., 2025): (1) Deriving general conclusions from specific cases, allowing it to cover and generalize to a wider range of applications, which aligns with the human learning process. (2) Adaptive adjustments help in uncertain and complex scenarios, where inductive reasoning may yield multiple plausible outcomes rather than a single unique solution.

## 3 Methodology

In this section, we will introduce three major approaches to enhance the inductive capabilities of LLMs: post-training enhancement (Section 3.1), test-time exploration (Section 3.2), and data augmentation (Section 3.3). It is worth noting that we not only summarize existing methods but also prospect a forward-looking review of potential future inductive approaches. For convenience, we treat the inputs of inductive reasoning as observations and refer to these outputs as rules. Boundaries and comparisons of the three methods can be found in the Appendix A.4.

### 3.1 Post-training Enhancement

Post-training enhancement refers to improving the inductive reasoning ability of LLMs during the

post-training stage (Lai et al., 2025; Wu, 2025), using algorithms such as supervised finetuning (SFT) and reinforcement learning (RL). This category of methods primarily focuses on constructing synthetic data (Long et al., 2024; Jiang et al., 2025) (Section 3.1.1) and developing new algorithms (Section 3.1.2). We illustrate them in Figure 3.

### 3.1.1 Synthetic Data

Synthetic data means artificially generated data that mimics the properties and patterns of real-world data (Bauer et al., 2024). Data plays a decisive role in LLM training. To address certain inherent limitations of natural data, such as being difficult to obtain or organize (Nadas et al., 2025), researchers often construct data manually to compensate for these shortcomings. LingR (Jiang et al., 2024a) builds a ‘linguistic rule instruction set’ for various LLMs, enabling them to learn step-by-step reasoning based on linguistic rules such as causality. ItD (Sun et al., 2024b) leverages the deductive abilities of LLMs to generate data and optimize inductive ability. The model’s capacity to learn general rules from a small number of samples is significantly enhanced. CodeSeq (Chen et al., 2025a) constructs SFT and RL training sets to ask LLMs to facilitate reasoning over number sequence general term formulas, thereby improving their inductive abilities. Other approaches (Wu et al., 2022; Aksu et al., 2023; Darm et al., 2023; Mosolova et al., 2025; Li et al., 2025b) establish similar induction-related training datasets for the models to learn from.

### 3.1.2 IRL-style Optimization

Reward models (RMs) (Zhong et al., 2025) are typically utilized to provide supervision signals for the RL process. However, for inductive reasoning, due to the non-uniqueness of answers and the uncertainty in the reasoning process, traditional RMs struggle to provide effective supervision. Therefore, Inverse RL (Arora and Doshi, 2021) (IRL), which needs to induce the latent reward functions, may serve as an alternative approach (Sun and van der Schaar, 2025). The Reinforcement Learning from Human Feedback (RLHF) (Kaufmann et al., 2024; Swamy et al., 2024) process of LLMs is essentially IRL, as it infers human preferences and the underlying reward function from human feedback. Therefore, designing an appropriate reward model in RLHF can enhance the inductive reasoning ability of LLMs (Jin et al., 2025). We can also employ Prompt-OIRL (Sun et al., 2024a) for

reference, which proposes an IRL-based method that reuses historical prompting trial-and-error experience to train a reward model to improve the model’s inductive exploration ability. Although works combining IRL and reasoning are still scarce, the approach of IRL—fitting the posterior distribution of the reward model from human or data signals (Cai et al., 2024; Krishna and Sahoo, 2024)—has strong extensibility and can thus be regarded as one of the important methods for the facilitation of inductive reasoning.

## 3.2 Test-time Exploration

The goal of inductive reasoning is to derive general rules from observations, which inevitably involves forming hypotheses during the reasoning process. The above post-training enhancement requires training LLMs, whereas the test-time exploration in this section is a hypothesis-based method that only works during the inference stage (Zhang et al., 2025c). Unlike end-to-end models (Kotary et al., 2021), the test-time exploration method prompts the frozen LLMs to form an inductive reasoning pipeline (He and Chen, 2025). We have LLMs generate candidate hypotheses for inductive problems, which can then undergo selection (Section 3.2.1), iteration (Section 3.2.2), or evolution (Section 3.2.3) operations to reach the optimal one. Detailed processes are shown in Figure 4.

### 3.2.1 Hypothesis Selection

Hypothesis selection refers to choosing, from the candidate hypotheses generated by LLMs, those that can cover the observations (Pazzani and Silverstein, 1990; Bun et al., 2019). Hypothesis Search (Wang et al., 2024b) let the LLMs generate multiple abstract hypotheses in natural language. Then, narrow down the hypothesis set through either the LLMs or minimal human filtering. The motivation of Mixture of Concepts (MoC) (Lee et al., 2025a) lies in the fact that hypotheses for inductive reasoning often produce semantic redundancy. Therefore, the proposed method simulates human inductive reasoning by first figuring out a list of semantically non-redundant concepts and then generating corresponding hypotheses based on each concept. Parfenova and Pfeffer (2025) proposes Ensemble Pipeline for Inductive Coding (EPIC) to address the issue of inconsistency in inductive encodings by using small LLMs to generate candidate encodings, filtering them through a moderator mechanism and similarity checks, and finally evaluating them with

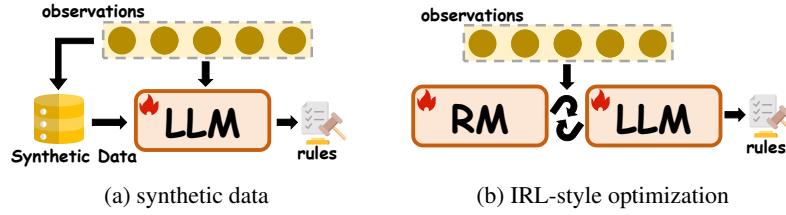


Figure 3: The demonstration of the post-training enhancement method for inductive reasoning.

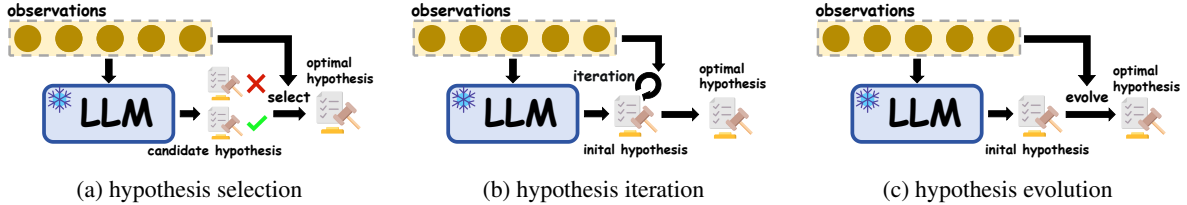


Figure 4: The demonstration of the test-time exploration method for inductive reasoning.

319 composite metrics.

### 320 3.2.2 Hypothesis Iteration

321 Hypothesis iteration means iterating over candidate  
 322 hypotheses until they satisfy all the observations  
 323 (Yom, 2015). Qiu et al. (2024) proposes a three-  
 324 step iterative hypothesis refinement method that  
 325 simulates the human inductive reasoning process:  
 326 generate multiple hypotheses from a few examples;  
 327 evaluate how many known instances each hypothe-  
 328 sis can cover; and have the LLMs further revise  
 329 the selected hypotheses based on feedback, iterat-  
 330 ing for several rounds until convergence. SSR (Li  
 331 et al., 2025a) iteratively optimizes by generating  
 332 diverse candidate rules and refining them based on  
 333 execution feedback. ARISE (M et al., 2025) also it-  
 334 erates over the inductive rules before using them to  
 335 train the model. IDEA framework (He et al., 2025)  
 336 resolves the shortcomings of LLMs in interactive  
 337 rule learning by simulating the human cycle of hy-  
 338 pothesis revision, thereby enhancing the model’s  
 339 dynamic learning capability.

### 340 3.2.3 Hypothesis Evolution

341 Unlike the iterative process, hypothesis evolution  
 342 expands, diversifies, or evolves the hypothesis  
 343 space by generating, filtering, and combining mul-  
 344 tiple hypotheses, forming hypotheses that better  
 345 capture complex patterns (Galkin, 2011; Gil et al.,  
 346 2017; Juretic, 2025). LLMs leverage contextual  
 347 and label information, along with prompts, to pro-  
 348 gressively guide the model in dynamically generat-  
 349 ing patterns during reasoning, without relying on  
 350 predefined rules (Dror et al., 2023). IncSchema (Li  
 351 et al., 2023a) gradually induces general patterns

352 by querying the LLMs in stages—first listing core  
 353 events, then expanding details, and finally verifying  
 354 relationships. HRI (Glanois et al., 2022) generates  
 355 inductive meta-rules and matches them with sam-  
 356 ples, thereby evolving into first-order logic rules.  
 357 PRIMO (Liu et al., 2024a) introduces a progres-  
 358 sive multi-stage open rule induction method for  
 359 deriving multi-hop rules, thereby capturing more  
 360 complex reasoning chains.

### 361 3.3 Data Augmentation

362 Data augmentation (Zhang et al., 2022) for LLMs  
 363 signifies enriching the model’s input with addi-  
 364 tional knowledge or structured signals, such as ex-  
 365 ternal facts and retrieved documents, to enhance  
 366 reasoning and improve output quality. We divide it  
 367 into three subcategories: human intervention (Sec-  
 368 tion 3.3.1), external knowledge (Section 3.3.2), and  
 369 structured signals (Section 3.3.3). Please refer to  
 370 Figure 5 about them.

#### 371 3.3.1 Human Intervention

372 Human intervention incorporates expert knowledge  
 373 or human-annotated information during inductive  
 374 reasoning. SS-VQ-VAE (Huang and Ji, 2020) relies  
 375 on a small amount of human-annotated information  
 376 to discover new patterns. Eyal et al. (2022) gener-  
 377 ates substitute words, then annotates the corpus and  
 378 trains static embeddings to enhance the model’s in-  
 379 ductive priors. Zhang et al. (2023a) utilizes GPT-3  
 380 to generate candidate patterns and enhances their  
 381 quality through human intervention, addressing the  
 382 issues of domain transfer and semantic consistency  
 383 in purely automated approaches. Some other stud-  
 384 ies (Edwards and Ji, 2023; Verhoeven et al., 2024)

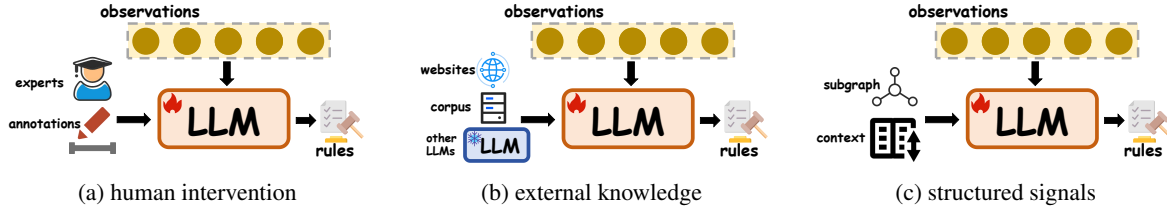


Figure 5: The demonstration of the data augmentation method for inductive reasoning.

emphasize inductive capabilities in low-annotation scenarios, which indirectly reflects the importance of expert knowledge.

### 3.3.2 External Knowledge

In this paper, we define external knowledge (Cao et al., 2021) to include web or document information, knowledge from other corpora, knowledge stored in LLM parameters (AlKhamissi et al., 2022), and so on. LLEGO (Liu et al., 2025) incorporates the semantic prior knowledge embedded in large LLMs into genetic programming operations to enhance generalization ability. The parameter knowledge of LLMs (Zhang et al., 2023b) and multimodal large models (Zhang et al., 2021b; Li et al., 2024a) can also serve as an important supplementary information source for inductive tasks. For example, some powerful LLMs are directly prompted to produce the inductive Chain-of-Thought (Chen et al., 2024c), inductive steps (Qian et al., 2023), and inductive rules (Ramji and Ramji, 2025) for the current task, providing additional assistance. Other types of knowledge, such as bilingual corpora (Shi et al., 2021; Kohli et al., 2024), social media content (Radhakrishnan et al., 2020), and commonsense knowledge (Ryu et al., 2022), can be used for inductive tasks in the same way.

### 3.3.3 Structured Signals

Structured information represents subgraphs or contextual information of LLMs (neighboring hidden states or embeddings), which provide local implicit signals and help LLMs to learn relevant inductive biases (Immer et al., 2022). Li et al. (2023b) optimizes the model’s output by retrieving nearest-neighbor embeddings as contextual examples. QARR (Xie et al., 2023) extracts an open subgraph for the query entity to inductively infer new entities. REST (Liu et al., 2024c) deploys rule-induced subgraphs to capture local semantic patterns, thereby enhancing the model’s generalization ability in inductive scenarios. GI-LUG (Kai et al., 2024) uses a syntactic parser to generate syn-

tax masks that guide the attention mechanism, and combine BPE embeddings with a hybrid loss function to optimize the induction process. Although this type of method is widely used in the PLM era, due to the same underlying principle, it can also play an important role for LLMs.

## 4 Evaluation

In this section, we will introduce current benchmarks for LLM inductive reasoning, some evaluation approaches, and the corresponding metrics.

### 4.1 Benchmarks

To evaluate the inductive reasoning capabilities of LLMs, the research community constructs a diverse set of benchmarks, as shown in Table 1, that comprehensively assess the models’ ability to generalize universal rules from concrete observations. It is noteworthy that the input formats of some tasks appear as paired samples or few-shot examples, often framed as analogy reasoning. As we claim in Appendix A.2, since analogical reasoning is a special form of inductive reasoning, we also regard benchmarks in this analogical form as benchmarks for inductive reasoning. The core task of these inductive benchmarks requires models to observe a small number of input examples (Observation Input), infer underlying patterns, and output the final rules (Induction Target).

As shown in Table 1, the data objects covered by these benchmarks span a wide range—from basic structures such as numbers, strings, and lists, to more complex forms like grids, logical formulas, and even natural language text.

Among them, benchmarks such as ARC (Chollet, 2019), List Functions (Rule, 2020), and SyGuS (Odena et al., 2021) focus on algorithmic or rule learning, requiring models to generate programs or operational rules that explain data transformations. What’s more, tasks like ILP (Evans and Grefenstette, 2018), GeoILP (Chen et al., 2025c), and ACRE (Zhang et al., 2021a) place greater emphasis on the induction of logical concepts and symbolic

Table 1: Some benchmarks for evaluating the inductive reasoning abilities of LLMs. We provide the atomic objects, names with their references, the input formations, the targets to be induced, and the number of test samples (approximate values). ‘.’ represents that it is the abbreviation of the benchmark name, while ‘\*’ indicates that the data are presented in the form of analogical reasoning. Further details about these benchmarks are in Appendix A.5.

Object	Benchmark Name	Observation Input	Induction Target	# Samples
symbol	ILP (Evans and Grefenstette, 2018)	pos. and neg. samples	a one-order logic rule	1,500
entity	SCAN (Lake and Baroni, 2018)	a state of entities	an action of the state	4,000
grid	ARC* (Chollet, 2019)	pairs of grids	a grid conversion rule	400
list	List Func.* (Rule, 2020)	pairs of number lists	a list operation rule	1,200
code	PROGES (Alet et al., 2021)	IO input/output	a program	270,000
string	SyGuS (Odena et al., 2021)	a pair of strings	a string-mapping program	8,000
entity	ACRE (Zhang et al., 2021a)	functions of entities	a ‘Blickets’ entity	6,000
text	Instruc. (Honovich et al., 2022)	two NL sentences	an instruction	2,400
number	Arith.* (Wu et al., 2024)	two numbers	the sum in certain base	1,000
symbol	Le/Ho. (Liu et al., 2024b)	pairs of triplets	an entailment rule	2,000
structure	NutFrame (Guo et al., 2024)	some frame information	conceptual structures	30,000
fact	DEER (Yang et al., 2024)	a pair of facts	a rule covers the facts	200
puzzle	RULEARN (He et al., 2025)	some puzzle scenarios	a puzzle rule	300
word	Crypto.* (Li et al., 2025a)	pairs of english words	an encrypted rule	300
symbol	GeoILP (Chen et al., 2025c)	pos. and neg. samples	a logic rule	250,000
string	In.Bench (Hua et al., 2025)	a pair of strings	a string-mapping rule	1,000
number	CodeSeq (Chen et al., 2025a)	a number sequence	the general term	200

rules. Particularly, Codeseq (Chen et al., 2025a) involves the computation of number sequence general terms, which represents a more advanced and complex form of inductive reasoning.

Overall, these benchmarks test the models’ pattern recognition ability and impose rigorous challenges on their higher-order cognitive skills. They not only examine how effectively LLMs can generalize from limited observations to underlying rules, but also serve as a foundation for driving further progress in enhancing such abilities.

## 4.2 Evaluation Approaches

In this section, we introduce the evaluation approaches for the inductive reasoning benchmarks mentioned above. We first present the existing evaluation strategies used in the benchmark papers (Section 4.2.1). Then, we derive a sandbox-based evaluation approach with a fine-grained observation coverage metric built upon it (Section 4.2.2).

### 4.2.1 Existing Evaluation Strategy

Most of the benchmarks in Section 4.1 and current works directly evaluate the consistency between answers generated by LLMs and the ground truth. Therefore, general metrics are employed, such as ACC, exact match, success rate, and so on. For example, ACRE (Zhang et al., 2021a) selects the most plausible “Blicket” from multiple options to

evaluate the accuracy of its selection. SCAN (Lake and Baroni, 2018) focuses on assessing whether the generated outputs exactly match the reference answers to indicate accuracy. SyGuS (Odena et al., 2021) requires finding a program that satisfies the string transformation rule, and the number of tasks in which the correct program is successfully identified is counted as the success rate.

### 4.2.2 Sandbox-based Evaluation

Considering that all inductive reasoning tasks share the same intrinsic mechanism: inferring general rules from specific observations. We can adopt a unified approach, namely sandbox unit test (Appendix A.7.1), to standardize the evaluation across all the inductive benchmarks mentioned above.

The sandbox unit test is a method where individual components or functions are tested in isolation to ensure they work as intended (Alhindi and Hallett, 2025). Each test runs in a controlled, independent environment, using specific input cases to verify the correctness of the component. This approach helps identify errors early and ensures that each part functions correctly before integration.

The sandbox unit test is originally used for code verification, at which time it is referred to as a code unit test (Gong et al., 2025). With the development of LLMs, it is also applied to evaluate various deductive reasoning and agent-related tasks of LLMs,

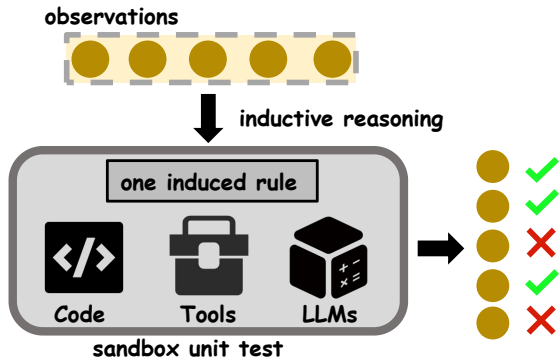


Figure 6: A demo of the sandbox unit test for inductive reasoning of LLMs.

such as InternBootCamp (Li et al., 2025c).

In inductive reasoning tasks, the sandbox unit test can also be deployed for evaluation. We present a demo in Figure 6. For an inductive rule generated by an LLM, it can be encapsulated as code (Gong et al., 2025), a tool (Qu et al., 2025), or written into a prompt to be provided to the LLM-as-a-Judge (Gu et al., 2024). Each observation can then be executed in a sandbox environment to determine whether it conforms to the current rule.

Based on this, we can derive a more fine-grained metric for LLM inductive reasoning: **observation coverage** (OC) (Appendix A.7.2), defined as the proportion of observations that pass the unit tests out of the total number of observations. In the example shown in Figure 6, this value is 0.6. Compared with the overall ACC or success rate of a task, OC provides a more fine-grained supervision signal at the observation level. This allows for a more precise reflection of the comprehensiveness of the model’s answer. With this metric, more informative feedback can be provided for subsequent rule refinement and hypothesis exploration.

## 5 Analysis

In this section, we present several prior exploratory tasks that offer theoretical analyses for inductive reasoning and inductive bias of LLMs (Kharitonov and Chaabouni, 2021; Papadimitriou and Jurafsky, 2023; Wilson and Frank, 2023).

**Inductive ability originates from induction heads.** Many studies (Si et al., 2023; Edelman et al., 2024; Chen et al., 2024d) show that the strong in-context learning (ICL) (Dong et al., 2023, 2024; Crosbie and Shutova, 2025) or example imitation (Honovich et al., 2023; Ye et al., 2025) ability of LLMs originates from induction heads. An induc-

tion head is an attention head (Edelman et al., 2022; Ren et al., 2024) that performs a match-and-copy operation, identifying and replicating relevant context tokens (Singh et al., 2024). Minegishi et al. (2025) finds that, in fact, induction heads are meta-learning an abstract inductive within the context.

**Model parameters, architecture, and data all help shape the inductive bias.** The parameters, model architecture, and training data (White and Cotterell, 2021; Merrill et al., 2021; Lovering et al., 2021; Levine et al., 2022; HaoChen and Ma, 2023; Movahedi et al., 2025) are key to inductive bias. Lippl and Lindsey (2024) explores the effects of different parameters on inductive bias under multi-data mixed training and single-task finetuning scenarios, and ultimately emphasizes the importance of task similarity in mixed training. Some studies (Cabannes et al., 2023; Aerni et al., 2023) also highlight the importance of data augmentation, even the noisy data. Further research (Zeno et al., 2025) demonstrates that the choice of minimum norm can also determine a model’s inductive generalization.

**Induction means simplicity.** Some early studies show that complex model architectures and data (Zietlow et al., 2021) can actually hinder inductive generalization. At the same time, for higher-order models, regularization can actually be detrimental to the formation of their inductive bias (Phuong and Lampert, 2021; Donhauser et al., 2022). Sometimes, simplicity is perfect for inductive reasoning (Goldblum et al., 2024). To enhance the inductive reasoning ability, finding simple inductive bias is of paramount importance. Simple and pure corpora often serve as the foundation for successful inductive reasoning (Mueller and Linzen, 2023).

## 6 Conclusion

This is the first survey of inductive reasoning for LLMs. The inductive mode is crucial for knowledge generalization and aligns better with human cognition. We categorize methods for improving inductive reasoning into three areas: post-training enhancement, test-time exploration, and data augmentation. We also summarize the current benchmarks and derive a unified sandbox-based evaluation approach. Finally, we offer some analyses regarding the source of inductive ability and how simple model architectures and data help with inductive tasks, providing a solid foundation for future research (Appendix A.8).

## 607 **Limitations**

608 This paper is a survey about the inductive reasoning  
609 abilities of Large Language Models. Due to space  
610 limitations, the main body of this survey is con-  
611 strained to fewer than eight pages, and therefore,  
612 many details are not included in the main text. We  
613 only present the most essential parts. Meanwhile,  
614 although inductive reasoning in LLMs attracts in-  
615 creasing attention in recent years, the number of  
616 related studies remains relatively limited, making it  
617 difficult to produce a large-scale, systematic survey  
618 (even extending to 100 pages) comparable to those  
619 in other areas.

## 620 **Ethics Statements**

621 This survey primarily organizes and summarizes ex-  
622 isting work on inductive reasoning in LLMs, with  
623 all relevant sources properly cited. Therefore, this  
624 paper does not raise any ethical or moral concerns.

## 625 **Acknowledgments**

626 We will finish this part in the camera-ready version.

## 627 **Use of AI Assistants**

628 We primarily use AI assistants to improve and en-  
629 rich our writing, especially by leveraging LLMs to  
630 help us write taxonomy in LaTeX.

## 631 **References**

632 Michael Aerni, Marco Milanta, Konstantin Donhauser,  
633 and Fanny Yang. 2023. [Strong inductive biases prov-  
634 ably prevent harmless interpolation](#). In *The Eleventh  
635 International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.  
636 OpenReview.net.  
637  
638 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui  
639 Zhang, and Wenpeng Yin. 2024. [Large language  
640 models for mathematical reasoning: Progresses and  
641 challenges](#). In *Proceedings of the 18th Conference of  
642 the European Chapter of the Association for Computa-  
643 tional Linguistics, EACL 2024: Student Research  
644 Workshop, St. Julian's, Malta, March 21-22, 2024*,  
645 pages 225–237. Association for Computational Lin-  
646 guistics.  
647  
648 Taha Aksu, Devamanyu Hazarika, Shikib Mehri,  
649 Seokhwan Kim, Dilek Hakkani-Tür, Yang Liu, and  
650 Mahdi Namazifar. 2023. [Cesar: Automatic induction  
651 of compositional instructions for multi-turn dialogs](#).  
*Preprint*, arXiv:2311.17376.  
652  
653 Ferran Alet, Javier Lopez-Contreras, James Koppel,  
654 Maxwell I. Nye, Armando Solar-Lezama, Tomás  
Lozano-Pérez, Leslie Pack Kaelbling, and Joshua B.

Tenenbaum. 2021. [A large-scale benchmark for few-  
shot program induction and synthesis](#). In *Proceed-  
ings of the 38th International Conference on Ma-  
chine Learning, ICML 2021, 18-24 July 2021, Vir-  
tual Event*, volume 139 of *Proceedings of Machine  
Learning Research*, pages 175–186. PMLR.

Maysara Alhindi and Joseph Hallett. 2025. [Playing in  
the sandbox: A study on the usability of seccomp](#).  
In *Twenty-First Symposium on Usable Privacy and  
Security, SOUPS 2025, Seattle, WA, USA, August  
10-12, 2025*, pages 225–240. USENIX Association.

Badr AlKhamissi, Millicent Li, Asli Celikyilmaz,  
Mona T. Diab, and Marjan Ghazvininejad. 2022. [A review on language models as knowledge bases](#).  
*CoRR*, abs/2204.06031.

Nadia Alshahwan, Jubin Chheda, Anastasia Finogenova,  
Beliz Gokkaya, Mark Harman, Inna Harper, Alexan-  
dru Marginean, Shubho Sengupta, and Eddy Wang.  
2024. [Automated unit test improvement using large  
language models at meta](#). *Companion Proceedings  
of the 32nd ACM International Conference on the  
Foundations of Software Engineering*.

Ohad Amosy, Tomer Volk, Eilam Shapira, Eyal  
Ben-David, Roi Reichart, and Gal Chechik. 2024. [Text2model: Text-based model induction for zero-  
shot image classification](#). In *Findings of the Associ-  
ation for Computational Linguistics: EMNLP 2024,  
Miami, Florida, USA, November 12-16, 2024*, pages  
155–172. Association for Computational Linguistics.

Saurabh Arora and Prashant Doshi. 2021. A survey of  
inverse reinforcement learning: Challenges, methods  
and progress. *Artificial Intelligence*, 297:103500.

W. Brian Arthur. 1994. [Inductive reasoning and  
bounded rationality](#). *The American Economic Re-  
view*, 84:406–411.

Yuyang Bai, Shangbin Feng, Vidhisha Balachandran,  
Zhaoxuan Tan, Shiqi Lou, Tianxing He, and Yulia  
Tsvetkov. 2024. [Kgquiz: Evaluating the generaliza-  
tion of encoded knowledge in large language models](#).  
In *Proceedings of the ACM on Web Conference 2024,  
WWW 2024, Singapore, May 13-17, 2024*, pages  
2226–2237. ACM.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,  
and Pascale Fung. 2023. [A multitask, multilingual,  
multimodal evaluation of chatgpt on reasoning, hal-  
lucination, and interactivity](#). *ArXiv*, abs/2302.04023.

Howard Barnum. 2012. [The beginning of infinity: Ex-  
planations that transform the world](#). *Physics Today*,  
65:48–50.

André Bauer, Simon Trapp, Michael Stenger, Robert  
Leppich, Samuel Kounev, Mark Leznik, Kyle Chard,  
and Ian Foster. 2024. [Comprehensive exploration of  
synthetic data generation: A survey](#). *arXiv preprint  
arXiv:2401.02524*.

711	Jonathan Baxter. 2000. <a href="#">A model of inductive bias learning</a> . <i>ArXiv</i> , abs/1106.0245.		
712			
713	Mark Bun, Gautam Kamath, Thomas Steinke, and Steven Z Wu. 2019. <a href="#">Private hypothesis selection</a> . <i>Advances in Neural Information Processing Systems</i> , 32.		
714			
715			
716			
717	Vivien Cabannes, Bobak Toussi Kiani, Randall Balestrieri, Yann LeCun, and Alberto Bietti. 2023. <a href="#">The SSL interplay: Augmentations, inductive bias, and generalization</a> . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 3252–3298. PMLR.		
718			
719			
720			
721			
722			
723			
724			
725	Xuan Cai, Xuesong Bai, Zhiyong Cui, Danmu Xie, Daocheng Fu, Haiyang Yu, and Yilong Ren. 2025. <a href="#">Text2scenario: Text-driven scenario generation for autonomous driving test</a> . <i>Preprint</i> , arXiv:2503.02911.		
726			
727			
728			
729			
730	Yuang Cai, Yuyu Yuan, Jinsheng Shi, and Qinrong Lin. 2024. <a href="#">Approximated variational bayesian inverse reinforcement learning for large language model alignment</a> . <i>Preprint</i> , arXiv:2411.09341.		
731			
732			
733			
734	Pengfei Cao, Xinyu Zuo, Yubo Chen, Kang Liu, Jun Zhao, Yuguang Chen, and Weihua Peng. 2021. <a href="#">Knowledge-enriched event causality identification via latent structure induction networks</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021</i> , pages 4862–4872. Association for Computational Linguistics.		
735			
736			
737			
738			
739			
740			
741			
742			
743			
744	Colin L. Carter and Howard J. Hamilton. 1998. <a href="#">Efficient attribute-oriented generalization for knowledge discovery from large databases</a> . <i>IEEE Trans. Knowl. Data Eng.</i> , 10(2):193–208.		
745			
746			
747			
748	Rich Caruana. 1993. <a href="#">Multitask learning: A knowledge-based source of inductive bias</a> . In <i>International Conference on Machine Learning</i> .		
749			
750			
751	Bowen Chen, Rune Sætre, and Yusuke Miyao. 2024a. <a href="#">A comprehensive evaluation of inductive reasoning capabilities and problem solving in large language models</a> . In <i>Findings of the Association for Computational Linguistics: EACL 2024, St. Julian's, Malta, March 17-22, 2024</i> , pages 323–339. Association for Computational Linguistics.		
752			
753			
754			
755			
756			
757			
758	Kedi Chen, Qin Chen, Jie Zhou, Yishen He, and Liang He. 2024b. <a href="#">Dialhalu: A dialogue-level hallucination evaluation benchmark for large language models</a> . <i>CoRR</i> , abs/2403.00896.		
759			
760			
761			
762	Kedi Chen, Zhikai Lei, Fan Zhang, Yinqi Zhang, Qin Chen, Jie Zhou, Liang He, Qipeng Guo, Kai Chen, and Wei Zhang. 2025a. <a href="#">Code-driven inductive synthesis: Enhancing reasoning abilities of large language models with sequences</a> . <i>Preprint</i> , arXiv:2503.13109.		
763			
764			
765			
766			
767			
	Po-Chun Chen, Sheng-Lun Wei, Hen-Hsen Huang, and Hsin-Hsi Chen. 2024c. <a href="#">Induct-learn: Short phrase prompting with instruction induction</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 5204–5231. Association for Computational Linguistics.		768
			769
			770
			771
			772
			773
			774
	Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025b. <a href="#">Towards reasoning era: A survey of long chain-of-thought for reasoning large language models</a> . <i>Preprint</i> , arXiv:2503.09567.		775
			776
			777
			778
			779
			780
	Si Chen, Richong Zhang, and Xu Zhang. 2025c. <a href="#">Geoilp: A synthetic dataset to guide large-scale rule induction</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.		781
			782
			783
			784
			785
	Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. 2024d. <a href="#">Unveiling induction heads: Provable training dynamics and feature learning in transformers</a> . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .		786
			787
			788
			789
			790
			791
			792
	Ziyang Chen, Dongfang Li, Xiang Zhao, Baotian Hu, and Min Zhang. 2024e. <a href="#">Temporal knowledge question answering via abstract reasoning induction</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 4872–4889. Association for Computational Linguistics.		793
			794
			795
			796
			797
			798
			799
			800
	Kewei Cheng, Jingfeng Yang, Haoming Jiang, Zhengyang Wang, Binxuan Huang, Ruirui Li, Shiyang Li, Zheng Li, Yifan Gao, Xian Li, Bing Yin, and Yizhou Sun. 2024a. <a href="#">Inductive or deductive? rethinking the fundamental reasoning abilities of llms</a> . <i>CoRR</i> , abs/2408.00114.		801
			802
			803
			804
			805
			806
	Sitao Cheng, Liangming Pan, Xunjian Yin, Xinyi Wang, and William Yang Wang. 2024b. <a href="#">Understanding the interplay between parametric and contextual knowledge for large language models</a> . <i>Preprint</i> , arXiv:2410.08414.		807
			808
			809
			810
			811
	François Chollet. 2019. <a href="#">On the measure of intelligence</a> . <i>CoRR</i> , abs/1911.01547.		812
			813
	Zhendong Chu, Shen Wang, Jian Xie, Tinghui Zhu, Yibo Yan, Jinheng Ye, Aoxiao Zhong, Xuming Hu, Jing Liang, Philip S. Yu, and Qingsong Wen. 2025. <a href="#">LLM agents for education: Advances and applications</a> . <i>CoRR</i> , abs/2503.11733.		814
			815
			816
			817
			818
	Herbert H. Clark. 1969. <a href="#">Linguistic processes in deductive reasoning</a> . <i>Psychological Review</i> , 76:387–404.		819
			820
	Irving M. Copi, Carl Cohen, and Daniel E. Flage. 2004. <a href="#">Essentials of logic</a> .		821
			822



933	Micah Goldblum, Marc Anton Finzi, Keefer Rowan, and Andrew Gordon Wilson. 2024. <a href="#">Position: The no free lunch theorem, kolmogorov complexity, and the role of inductive biases in machine learning</a> . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	988	Or Honovich, Uri Shaham, Samuel R. Bowman, and Omer Levy. 2023. <a href="#">Instruction induction: From few examples to natural language task descriptions</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1935–1952. Association for Computational Linguistics.	998
934		990		999
935		991		992
936		992		993
937		993		994
938		994		995
939		995		
940	Jingzhi Gong, Vardan Voskanyan, Paul Brookes, Fan Wu, Wei Jie, Jie Xu, Rafail Giavrimis, Mike Basios, Leslie Kanthan, and Zheng Wang. 2025. <a href="#">Language models for code optimization: Survey, challenges and future directions</a> . <i>CoRR</i> , abs/2501.01277.	996	Wenyue Hua, Tyler Wong, Fei Sun, Liangming Pan, Adam Jardine, and William Yang Wang. 2025. <a href="#">Inductionbench: Llms fail in the simplest complexity class</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 26526–26546. Association for Computational Linguistics.	997
941		998		999
942		1000		1001
943		1001		1002
944		1002		1003
945	Anirudh Goyal and Yoshua Bengio. 2020. <a href="#">Inductive biases for deep learning of higher-level cognition</a> . <i>CoRR</i> , abs/2011.15091.	1003		
946		1004		1005
947		1005		1006
948	Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and Jian Guo. 2024. <a href="#">A survey on llm-as-a-judge</a> . <i>CoRR</i> , abs/2411.15594.	1006	Jie Huang and Kevin Chen-Chuan Chang. 2022. <a href="#">Towards reasoning in large language models: A survey</a> . <i>arXiv preprint arXiv:2212.10403</i> .	1007
949		1007		1008
950		1008		1009
951		1009		1010
952		1010		1011
953	Yu Gu and Yu Su. 2022. <a href="#">Arcaneqa: Dynamic program induction and contextualized encoding for knowledge base question answering</a> . In <i>Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022</i> , pages 1718–1731. International Committee on Computational Linguistics.	1011	Lifu Huang and Heng Ji. 2020. <a href="#">Semi-supervised new event type induction and event detection</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 718–724, Online. Association for Computational Linguistics.	1012
954		1012		1013
955		1013		1014
956		1014		1015
957		1015		1016
958		1016		1017
959		1017		1018
960	Shaoru Guo, Yubo Chen, Kang Liu, Ru Li, and Jun Zhao. 2024. <a href="#">Nutframe: Frame-based conceptual structure induction with llms</a> . In <i>LREC/COLING</i> , pages 12330–12335.	1018	Alexander Immer, Lucas Torroba Hennigen, Vincent Fortuin, and Ryan Cotterell. 2022. <a href="#">Probing as quantifying inductive bias</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 1839–1851. Association for Computational Linguistics.	1019
961		1019		1020
962		1020		1021
963		1021		
964	Simon Jerome Han, Keith J. Ransom, Andrew Perfors, and Charles Kemp. 2024. <a href="#">Inductive reasoning in humans and large language models</a> . <i>Cogn. Syst. Res.</i> , 83:101155.	1021	Tim Ingold. 2021. <a href="#">The perception of the environment: essays on livelihood, dwelling and skill</a> .	1022
965		1022		1023
966		1023		
967		1024		1025
968	Jeff Z. HaoChen and Tengyu Ma. 2023. <a href="#">A theoretical study of inductive biases in contrastive learning</a> . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1025	InternTeam. 2025. <a href="#">Intern-s1: A scientific multimodal foundation model</a> . <i>CoRR</i> , abs/2508.15763.	1026
969		1026		1027
970		1027		1028
971		1028		1029
972		1029		1030
973	Kaiyu He and Zhiyu Chen. 2025. <a href="#">From reasoning to learning: A survey on hypothesis discovery and rule learning with large language models</a> . <i>Preprint</i> , arXiv:2505.21935.	1030	Shuoran Jiang, Qingcai Chen, Yang Xiang, Youcheng Pan, and Yukang Lin. 2024a. <a href="#">Linguistic rule induction improves adversarial and OOD robustness in large language models</a> . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 10565–10577, Torino, Italia. ELRA and ICCL.	1031
974		1031		1032
975		1032		1033
976		1033		1034
977	Kaiyu He, Mian Zhang, Shuo Yan, Peilin Wu, and Zhiyu Zoey Chen. 2025. <a href="#">Idea: Enhancing the rule learning ability of large language model agent through induction, deduction, and abduction</a> . <i>Preprint</i> , arXiv:2408.10455.	1034	Xue Jiang, Yihong Dong, Lecheng Wang, Zheng Fang, Qiwei Shang, Ge Li, Zhi Jin, and Wenpin Jiao. 2024b. <a href="#">Self-planning code generation with large language models</a> . <i>ACM Trans. Softw. Eng. Methodol.</i> , 33(7):182:1–182:30.	1035
978		1035		1036
979		1036		
980		1037		1038
981		1038		1039
982	Evan Heit. 2000. <a href="#">Properties of inductive reasoning</a> . <i>Psychonomic Bulletin &amp; Review</i> , 7:569–592.	1039	Yanru Jiang, Siyu Liang, and Junwon Choi. 2025. <a href="#">Synthetic survey data generation and evaluation</a> . In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining, V.1, KDD 2025, Toronto, ON, Canada, August 3-7, 2025</i> , pages 2292–2302. ACM.	1040
983		1040		1041
984	Or Honovich, Uri Shaham, Samuel R. Bowman, and Omer Levy. 2022. <a href="#">Instruction induction: From few examples to natural language task descriptions</a> . <i>Preprint</i> , arXiv:2205.10782.	1041		1042
985		1042		
986				
987				

1043	Can Jin, Yang Zhou, Qixin Zhang, Hongwu Peng,	James Kotary, Ferdinando Fioretto, Pascal Van Hen-	1097
1044	Di Zhang, Marco Pavone, Ligong Han, Zhang-Wei	tenryck, and Bryan Wilder. 2021. <a href="#">End-to-end con-</a>	1098
1045	Hong, Tong Che, and Dimitris N. Metaxas. 2025.	<a href="#">strained optimization learning: A survey</a> . In <i>Proceed-</i>	1099
1046	<a href="#">Your reward function for rl is your best prm for</a>	<i>ings of the Thirtieth International Joint Conference</i>	1100
1047	<a href="#">search: Unifying rl and search-based tts</a> . <i>Preprint</i> ,	<i>on Artificial Intelligence, IJCAI 2021, Virtual Event /</i>	1101
1048	arXiv:2508.14313.	<i>Montreal, Canada, 19-27 August 2021</i> , pages 4475–	1102
		4482. ijcai.org.	1103
1049	Davor Juretic. 2025. <a href="#">Exploring the evolution-coupling</a>	Shambhavi Krishna and Aishwarya Sahoo. 2024. <a href="#">Solv-</a>	1104
1050	<a href="#">hypothesis: Do enzymes' performance gains corre-</a>	<a href="#">ing the inverse alignment problem for efficient rlhf</a> .	1105
1051	<a href="#">late with increased dissipation?</a> <i>Entropy</i> , 27(4):365.	<i>Preprint</i> , arXiv:2412.10529.	1106
1052	Jushi Kai, Shengyuan Hou, Yusheng Huang, and	Hanyu Lai, Xiao Liu, Junjie Gao, Jiale Cheng, Zehan	1107
1053	Zhouhan Lin. 2024. <a href="#">Leveraging grammar induction</a>	Qi, Yifan Xu, Shuntian Yao, Dan Zhang, Jinhua Du,	1108
1054	<a href="#">for language understanding and generation</a> . In <i>Find-</i>	Zhenyu Hou, Xin Lv, Minlie Huang, Yuxiao Dong,	1109
1055	<i>ings of the Association for Computational Linguis-</i>	and Jie Tang. 2025. <a href="#">A survey of post-training scaling</a>	1110
1056	<i>tics: EMNLP 2024, Miami, Florida, USA, November</i>	<a href="#">in large language models</a> . In <i>Proceedings of the 63rd</i>	1111
1057	<i>12-16, 2024</i> , pages 4501–4513. Association for Com-	<i>Annual Meeting of the Association for Computational</i>	1112
1058	putational Linguistics.	<i>Linguistics (Volume 1: Long Papers), ACL 2025, Vi-</i>	1113
1059	Charles W. Kalish and Jordan Thevenow-Harrison. 2014.	<i>enna, Austria, July 27 - August 1, 2025</i> , pages 2771–	1114
1060	<a href="#">Descriptive and inferential problems of induction</a> .	2791. Association for Computational Linguistics.	1115
1061	Katikapalli Subramanyam Kalyan, Ajit Rajasekharan,	Viet Dac Lai, Hieu Man, Linh Ngo Van, Franck Dernon-	1116
1062	and Sivanesan Sangeetha. 2021. <a href="#">AMMUS : A survey</a>	court, and Thien Huu Nguyen. 2022. <a href="#">Multilingual</a>	1117
1063	<a href="#">of transformer-based pretrained models in natural</a>	<a href="#">subevent relation extraction: A novel dataset and</a>	1118
1064	<a href="#">language processing</a> . <i>CoRR</i> , abs/2108.05542.	<a href="#">structure induction method</a> . In <i>Findings of the Asso-</i>	1119
1065	Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke	<i>ciation for Computational Linguistics: EMNLP 2022,</i>	1120
1066	Hüllermeier. 2024. <a href="#">A survey of reinforcement learn-</a>	<i>Abu Dhabi, United Arab Emirates, December 7-11,</i>	1121
1067	<a href="#">ing from human feedback</a> .	<i>2022</i> , pages 5559–5570. Association for Computa-	1122
1068	Ruhma Khan, Sumit Gulwani, Vu Le, Arjun Radhakr-	tional Linguistics.	1123
1069	ishna, Ashish Tiwari, and Gust Verbruggen. 2025.	Brenden M. Lake and Marco Baroni. 2018. <a href="#">General-</a>	1124
1070	<a href="#">Llm-guided compositional program synthesis</a> . <i>CoRR</i> ,	<a href="#">ization without systematicity: On the compositional</a>	1125
1071	abs/2503.15540.	<a href="#">skills of sequence-to-sequence recurrent networks</a> .	1126
1072	Eugene Kharitonov and Rahma Chaabouni. 2021. <a href="#">What</a>	<i>Preprint</i> , arXiv:1711.00350.	1127
1073	<a href="#">they do when in doubt: a study of inductive biases in</a>	Kang-il Lee, Hyukhun Koh, Dongryeol Lee, Seunghyun	1128
1074	<a href="#">seq2seq learners</a> . In <i>9th International Conference on</i>	Yoon, Minsung Kim, and Kyomin Jung. 2025a. <a href="#">Gen-</a>	1129
1075	<i>Learning Representations, ICLR 2021, Virtual Event,</i>	<a href="#">erating diverse hypotheses for inductive reasoning</a> .	1130
1076	<i>Austria, May 3-7, 2021</i> . OpenReview.net.	In <i>Proceedings of the 2025 Conference of the Na-</i>	1131
1077	Jeonghwan Kim, Giwon Hong, Sung-Hyon Myaeng,	<i>tions of the Americas Chapter of the Association for</i>	1132
1078	and Joyce Jiyoung Whang. 2023. <a href="#">Fineprompt: Un-</a>	<i>Computational Linguistics: Human Language Tech-</i>	1133
1079	<a href="#">veiling the role of finetuned inductive bias on com-</a>	<i>nologies, NAACL 2025 - Volume 1: Long Papers,</i>	1134
1080	<a href="#">positional reasoning in GPT-4</a> . In <i>Findings of the</i>	<i>Albuquerque, New Mexico, USA, April 29 - May 4,</i>	1135
1081	<i>Association for Computational Linguistics: EMNLP</i>	<i>2025</i> , pages 8461–8474. Association for Computa-	1136
1082	<i>2023, Singapore, December 6-10, 2023</i> , pages 3763–	tional Linguistics.	1137
1083	3775. Association for Computational Linguistics.	Kangil Lee, Jahyun Koo, Seunghyun Yoon, Minbeom	1138
1084	Taeuk Kim, Jihun Choi, Daniel Edmiston, and Sang-goo	Kim, Hyukhun Koh, Dongryeol Lee, and Kyomin	1139
1085	Lee. 2020. <a href="#">Are pre-trained language models aware</a>	Jung. 2025b. <a href="#">Program synthesis via test-time trans-</a>	1140
1086	<a href="#">of phrases? simple but strong baselines for gram-</a>	<a href="#">duction</a> . <i>CoRR</i> , abs/2509.17393.	1141
1087	<a href="#">mar induction</a> . In <i>8th International Conference on</i>	Yoav Levine, Noam Wies, Daniel Jannai, Dan Navon,	1142
1088	<i>Learning Representations, ICLR 2020, Addis Ababa,</i>	Yedid Hoshen, and Amnon Shashua. 2022. <a href="#">The in-</a>	1143
1089	<i>Ethiopia, April 26-30, 2020</i> . OpenReview.net.	<a href="#">ductive bias of in-context learning: Rethinking pre-</a>	1144
1090	Harsh Kohli, Helian Feng, Nicholas Dronen, Calvin	<a href="#">training example design</a> . In <i>The Tenth International</i>	1145
1091	McCarter, Sina Moeini, and Ali Kebarighotbi. 2024.	<i>Conference on Learning Representations, ICLR 2022,</i>	1146
1092	<a href="#">How lexical is bilingual lexicon induction?</a> In <i>Find-</i>	<i>Virtual Event, April 25-29, 2022</i> . OpenReview.net.	1147
1093	<i>ings of the Association for Computational Linguis-</i>	Martha Lewis and Melanie Mitchell. 2024a. <a href="#">Evaluat-</a>	1148
1094	<i>tics: NAACL 2024, Mexico City, Mexico, June 16-21,</i>	<a href="#">ing the robustness of analogical reasoning in large</a>	1149
1095	<i>2024</i> , pages 4381–4386. Association for Computa-	<a href="#">language models</a> . <i>arXiv preprint arXiv:2411.14215</i> .	1150
1096	tional Linguistics.	Martha Lewis and Melanie Mitchell. 2024b. <a href="#">Using</a>	1151
		<a href="#">counterfactual tasks to evaluate the generality of ana-</a>	1152
		<a href="#">logical reasoning in large language models</a> . <i>arXiv</i>	1153
		<i>preprint arXiv:2402.08955</i> .	1154

1155	Boyi Li, Rodolfo Corona, Karttikeya Mangalam, Catherine Chen, Daniel Flaherty, Serge Belongie, Kilian Q. Weinberger, Jitendra Malik, Trevor Darrell, and Dan Klein. 2024a. <a href="#">Re-evaluating the need for multimodal signals in unsupervised grammar induction</a> . <i>Preprint</i> , arXiv:2212.10564.	<i>Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 11558–11573. Association for Computational Linguistics.	1212 1213 1214 1215
1161	Chunyang Li, Weiqi Wang, Tianshi Zheng, and Yangqiu Song. 2025a. <a href="#">Patterns over principles: The fragility of inductive reasoning in llms under noisy observations</a> . In <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 19608–19626. Association for Computational Linguistics.	Samuel Lippl and Jack W. Lindsey. 2024. <a href="#">Inductive biases of multi-task learning and finetuning: multiple regimes of feature reuse</a> . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	1216 1217 1218 1219 1220 1221 1222
1166	Jiachun Li, Pengfei Cao, Zhuoran Jin, Yubo Chen, Kang Liu, and Jun Zhao. 2025b. <a href="#">MIRAGE: evaluating and explaining inductive reasoning process in language models</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	Jianyu Liu, Sheng Bi, and Guilin Qi. 2024a. <a href="#">Primo: Progressive induction for multi-hop open rule generation</a> . <i>Preprint</i> , arXiv:2411.01205.	1223 1224 1225
1171	Peiji Li, Jiasheng Ye, Yongkang Chen, Yichuan Ma, Zijie Yu, Kedi Chen, Ganqu Cui, Haozhan Li, Jiacheng Chen, Chengqi Lyu, Wenwei Zhang, Linyang Li, Qipeng Guo, Dahua Lin, Bowen Zhou, and Kai Chen. 2025c. <a href="#">Internbootcamp technical report: Boosting LLM reasoning with verifiable task scaling</a> . <i>CoRR</i> , abs/2508.08636.	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. <a href="#">Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1226 1227 1228 1229 1230 1231 1232 1233
1181	Sha Li, Ruining Zhao, Manling Li, Heng Ji, Chris Callison-Burch, and Jiawei Han. 2023a. <a href="#">Open-domain hierarchical event schema induction by incremental prompting and verification</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5677–5697, Toronto, Canada. Association for Computational Linguistics.	Tennison Liu, Nicolas Huynh, and Mihaela van der Schaar. 2025. <a href="#">Decision tree induction through llms via semantically-aware evolution</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	1234 1235 1236 1237 1238 1239
1189	Siyue Li, Xiaofan Zhou, Zhizhong Wu, Yuiian Long, and Yanxin Shen. 2024b. <a href="#">Strategic deductive reasoning in large language models: A dual-agent approach</a> . In <i>2024 IEEE 6th International Conference on Power, Intelligent Computing and Systems (ICPICS)</i> , pages 834–839. IEEE.	Tianyang Liu, Tianyi Li, Liang Cheng, and Mark Steedman. 2024b. <a href="#">Explicit inductive inference using large language models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024</i> , pages 15779–15786. Association for Computational Linguistics.	1240 1241 1242 1243 1244 1245
1195	Wen-Ding Li and Kevin Ellis. 2024. <a href="#">Is programming by example solved by llms?</a> In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Tianyu Liu, Qitan Lv, Jie Wang, Shuling Yang, and Hanzhu Chen. 2024c. <a href="#">Learning rule-induced subgraph representations for inductive relation prediction</a> . <i>CoRR</i> , abs/2408.07088.	1246 1247 1248 1249
1201	Yaoyiran Li, Anna Korhonen, and Ivan Vulic. 2023b. <a href="#">On bilingual lexicon induction with large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 9577–9599. Association for Computational Linguistics.	Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. 2024. <a href="#">On llms-driven synthetic data generation, curation, and evaluation: A survey</a> . <i>CoRR</i> , abs/2406.15126.	1250 1251 1252 1253
1208	Matthias Lindemann, Alexander Koller, and Ivan Titov. 2024. <a href="#">Strengthening structural inductive biases by pre-training to perform syntactic transformations</a> . In <i>Proceedings of the 2024 Conference on Empirical</i>	Charles Lovering, Rohan Jha, Tal Linzen, and Ellie Pavlick. 2021. <a href="#">Predicting inductive biases of pre-trained models</a> . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	1254 1255 1256 1257 1258
1211		Zhou Lu. 2024. <a href="#">When is inductive inference possible?</a> In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	1259 1260 1261 1262 1263
		Zimu Lu, Aojun Zhou, Ke Wang, Houxing Ren, Weikang Shi, Junting Pan, Mingjie Zhan, and Hongsheng Li. 2024. <a href="#">Mathcoder2: Better math reasoning from continued pretraining on model-translated mathematical code</a> . <i>CoRR</i> , abs/2410.08196.	1264 1265 1266 1267 1268



1382	Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Bäck. 2024. Reasoning with large language models, a survey. <i>CoRR</i> , abs/2407.11511.	1439
1383		1440
1384		1441
1385		1442
1386	Foster J. Provost and Bruce G. Buchanan. 1995. Inductive policy: The pragmatics of bias selection. <i>Machine Learning</i> , 20:35–61.	1443
1387		1444
1388		
1389	Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. 2023. Limitations of language models in arithmetic and symbolic induction. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2023, Toronto, Canada, July 9-14, 2023, pages 9285–9298. Association for Computational Linguistics.	1445
1390		1446
1391		1447
1392		
1393		1448
1394		1449
1395		1450
1396		1451
1397	Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin Wang, Yoon Kim, Yejin Choi, Nouha Dziri, and Xiang Ren. 2024. Phenomenal yet puzzling: Testing inductive reasoning capabilities of language models with hypothesis refinement. <i>Preprint</i> , arXiv:2310.08559.	1452
1398		1453
1399		1454
1400		1455
1401		
1402		1456
1403		1457
1404	Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2025. Tool learning with large language models: a survey. <i>Frontiers Comput. Sci.</i> , 19(8):198343.	1458
1405		1459
1406		
1407		1460
1408	QwenTeam. 2025. Qwen3 technical report. <i>CoRR</i> , abs/2505.09388.	1461
1409		
1410	Karthik Radhakrishnan, Tushar Kanakagiri, Sharanya Chakravarthy, and Vidhisha Balachandran. 2020. "a little birdie told me ..." - inductive biases for rumour stance detection on social media. In <i>Proceedings of the Sixth Workshop on Noisy User-generated Text, W-NUT@EMNLP 2020 Online, November 19, 2020</i> , pages 244–248. Association for Computational Linguistics.	1462
1411		1463
1412		1464
1413		1465
1414		1466
1415		1467
1416		1468
1417		1469
1418	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Preprint</i> , arXiv:1910.10683.	1470
1419		
1420		1471
1421		1472
1422		1473
1423	Raghav Ramji and Keshav Ramji. 2025. Inductive linguistic reasoning with large language models. In <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 22783–22810. Association for Computational Linguistics.	1474
1424		1475
1425		1476
1426		1477
1427		1478
1428		1479
1429	Chandan K. Reddy and Parshin Shojaee. 2025. Towards scientific discovery with generative AI: progress, opportunities, and challenges. In <i>AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA</i> , pages 28601–28609. AAAI Press.	1480
1430		1481
1431		
1432		1482
1433		1483
1434		1484
1435	Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. 2024. Identifying semantic induction heads to understand in-context learning. In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 6916–6932. Association for Computational Linguistics.	1485
1436		1486
1437		1487
1438		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500





1723	Chi Zhang, Baoxiong Jia, Mark Edmonds, Song-Chun Zhu, and Yixin Zhu. 2021a. <a href="#">ACRE: abstract causal reasoning beyond covariation</a> . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021</i> , pages 10643–10653. Computer Vision Foundation / IEEE.	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2025. <a href="#">A survey of large language models</a> . <i>Preprint</i> , arXiv:2303.18223.	1780
1724			1781
1725			1782
1726			1783
1727			1784
1728			
1729	Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, Yu Fu, Xingtai Lv, Yuchen Zhang, Sihang Zeng, Shang Qu, Haozhan Li, Shijie Wang, Yuru Wang, Xinwei Long, and 20 others. 2025a. <a href="#">A survey of reinforcement learning for large reasoning models</a> . <i>Preprint</i> , arXiv:2509.08827.	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. <a href="#">A survey of large language models</a> . <i>arXiv preprint arXiv:2303.18223</i> , 1(2).	1785
1730			1786
1731			1787
1732			1788
1733			1789
1734		Jialun Zhong, Wei Shen, Yanzeng Li, Songyang Gao, Hua Lu, Yicheng Chen, Yang Zhang, Wei Zhou, Jinjie Gu, and Lei Zou. 2025. A comprehensive survey of reward models: Taxonomy, applications, challenges, and future. <i>arXiv preprint arXiv:2504.12328</i> .	1790
1735			1791
1736			1792
1737			1793
1738			1794
1739			
1740			
1741			1795
1742			1796
1743			1797
1744			1798
1745			1799
1746			1800
1747			
1748			
1749			
1750			
1751			
1752			
1753			
1754			
1755			
1756			
1757			
1758			
1759			
1760			
1761			
1762			
1763			
1764			
1765			
1766			
1767			
1768			
1769			
1770			
1771			
1772			
1773			
1774			
1775			
1776			
1777			
1778			
1779			
1780			
1801			
1802			
1803			
1804			
1805			
1806			
1807			
1808			
1809			

## A Appendix

### A.1 Inductive Reasoning and Inductive Bias

The differences between inductive reasoning and inductive bias can be discussed in the following three aspects:

**Different concepts** Inductive reasoning is a type of reasoning task whose input consists of some observations and whose output is a general rule or principle that is consistent with these observations. In contrast, inductive bias is a concept that represents the prior knowledge or inclination within a machine learning system (Goyal and Bengio, 2020). It is the preferences or constraints that are ‘pre-specified’ in the model before learning. The process of learning such bias is called meta learning (Vanschoren, 2018) in the machine learning field.

**Different scopes of application** All machine learning and deep learning tasks (including text, vision, video, and other tasks) inherently require the model or algorithm to be injected with the inductive bias appropriate for the specific task. Inductive reasoning, however, refers to one of the tasks.

**Context of discussion** The goal of this inductive reasoning task is to identify rules among observations within a single task. This paper primarily discusses the inductive reasoning task within the context of LLMs. Improving the inductive reasoning ability of LLMs essentially amounts to discovering the inductive bias underlying inductive reasoning. By doing so, the scope of the discussion can be narrowed, allowing the survey to be more thorough and well-developed.

### A.2 Inductive Reasoning and Other Reasoning Modes

#### A.2.1 Inductive Reasoning and Deductive Reasoning

We provide an example for each in Figure 7. As two major modes of reasoning, inductive reasoning (Clark, 1969; Rips, 1994; Bang et al., 2023) differs from deductive reasoning in many aspects. Please refer to Table 2 for more.

#### A.2.2 Inductive Reasoning and Analogical Reasoning

As stated in the main text, inductive reasoning is a process that goes from the particular to the general. In contrast, analogical reasoning is a process from the particular to the particular, which can be

**Inductive Reasoning**  
Please provide the general term for the following number sequence:  
-1, 1, -1, 1, -1, 1, -1, 1, -1, 1, ...

answer1  
the general term is:  $(-1)^n$  ( $n \geq 1$ )      correct answer  
answer2  
the general term is:  $\cos(n \cdot \pi)$  ( $n \geq 1$ )      correct answer

**Deductive Reasoning**  
Please provide the answer to the following math problem:  
Given a square with a side length of 5 cm, what is the length of its diagonal?

answer1  
len(diagonal) is  $\sqrt{2}$  times len(side), so the answer:  $5 \cdot \sqrt{2}$ .      correct answer  
answer2  
the length of its diagonal is 10.      wrong answer

Figure 7: An example is given for both reasoning modes. In inductive reasoning, there may be multiple correct answers consistent with the existing observations, while in deductive reasoning, a precise logical reasoning process can help arrive at the only correct answer.

Table 2: The differences between inductive reasoning and deductive reasoning.

	Inductive Reasoning	Deductive Reasoning
<b>Mode</b>	particular-to-general	general-to-particular
<b>Target</b>	probabilistic conclusion	precise answer
<b>Reliability</b>	flexible and generative	premises-relied
<b>Informative</b>	extended information	no new information
<b>Application</b>	hypothetical inference	rigorous proof

regarded as a special form of inductive reasoning. For example, consider an inductive reasoning task of deriving the general term formula of a number sequence: the input is a number sequence, and the output is its general term formula. Analogical reasoning, on the other hand, takes a number sequence as input and outputs its next term. In short, analogical reasoning is a form of reasoning that involves imitation based on observations. It compares and transfers similarities from existing observations to infer and generate possible next items or outcomes (Lewis and Mitchell, 2024a,b; Musker et al., 2024).

In the field of LLMs, the commonly studied ICL (Dong et al., 2023, 2024) can be regarded as a form of analogical reasoning. For ease of discussion and to clearly define the scope, a large body of research on ICL is not discussed in this paper.

### A.3 More Synthesis and Comparative Analysis on NLP Applications

The impressive applications we cited, spanning from core syntactic and semantic parsing to complex multimodal tasks, clearly demonstrate the versatility of inductive reasoning in NLP.

**Empirical effectiveness** often hinges on the

method’s ability to seamlessly integrate domain-specific structural knowledge as inductive bias — for instance, achieving high precision in information extraction usually involves explicit rule-based or graph-structured inductive biases, while improvements in dialogue systems frequently rely on incorporating conversational history as an inductive structure.

Approaches that embed strong, hand-crafted inductive bias (e.g., explicit grammar rules for parsing) tend to offer **superior interpretability and data efficiency** but suffer from lower generalizability across domains, whereas soft or learned bias (e.g., the knowledge encoded in the model’s parameters) is **more flexible but potentially less transparent**. How to achieve a **trade-off** between the two is currently a key focus for applying inductive reasoning to downstream NLP tasks.

**The major trend** of inductive reasoning in NLP is a move toward methods that automate the discovery of inductive patterns, such as applying self-supervision or meta-learning to identify bias suitable for multiple tasks, marking a shift from prescriptive to descriptive induction (Kalish and Thevenow-Harrison, 2014) in method design.

This comparative analysis underscores the necessity of selecting an inductive reasoning method that carefully balances the need for robustness, interpretability, and generalization capacity for the target application scenario.

## A.4 The Boundaries and Comparisons among the Methods

### A.4.1 The Boundaries

Our classification minimizes overlap among the three methods conceptually as much as possible.

(1) **Post-training enhancement** directly trains LLMs using task-specific data. Synthetic data in this category is often generated from scratch to create entirely new training samples. (2) **Test-time exploration** is a train-free method based on inductive hypotheses. (3) **Data augmentation** can guide inductive tasks by leveraging annotated data and knowledge from outside the task domain manually. It can keep the model frozen, or even be used to train with fine-grained embedding or graph information. In conclusion, it modifies and expands existing training samples.

We believe that modern LLM techniques are often multifunctional, and some of the overlapping parts can be seen as emerging trends or directions

for future research. For example, the construction of synthetic data inevitably incorporates prior knowledge from human experts.

### A.4.2 The Comparisons

(1) **Post-training enhancement** is suitable for scenarios where sufficient data is available, and there is a need to directly enhance a single model’s inherent inductive capabilities. *Advantages*: It can directly enhance the inductive reasoning ability of a specific LLM and is easy and quick to train. *Disadvantages*: It is only applicable to scenarios involving a single model, and cannot be easily used to solve complex problems that require cooperation among agent systems.

(2) **Test-time exploration** can be applied in situations with limited data or when the model cannot be fine-tuned, using hypothesis-driven methods to perform inductive reasoning tasks. *Advantages*: No model training is required, as it can leverage the capabilities of powerful closed-source LLMs. *Disadvantages*: It needs time and monetary costs, and multiple refinements do not necessarily guarantee an optimal solution.

(3) **Data augmentation** is more appropriate for scenarios where the model involves cross-domain tasks or requires the use of external knowledge (e.g., cross-domain agent systems with tool usage and human-crafted rules), or for training auxiliary modules such as embeddings or graphs, thereby improving the model’s generalization and transferability to new patterns. *Advantages*: It can leverage information beyond the LLM, including even parameter-level information, to provide assistance. *Disadvantages*: Its working mode is similar to that of an agent system, giving it potential to handle complex real-world problems. However, it relies on external data and tools, and the overall system performance depends on the fundamental capabilities of each component.

### A.4.3 Why Synthetic Data is Widely Used

Synthetic data accounts for a significant portion of many well-known LLM training corpora (DeepSeek-AI, 2024; QwenTeam, 2025; InternTeam, 2025) and has the following advantages compared with real data.

(1) **Coverage of rare cases**: Synthetic data can include samples that are uncommon in real-world scenarios. For example, in a batch of real data, certain inductive patterns may appear infrequently, and directly training on such data could introduce

1981	bias into the model. (2) <b>Support for data-scarce or privacy-sensitive domains:</b> In fields where data is limited or highly sensitive, synthetic data can be generated at a large scale, enabling the training of more generalizable models. (3) <b>High efficiency:</b> Synthetic data can be generated very flexibly and at a large scale, which makes model training and debugging more efficient. (4) <b>Sometimes better performance:</b> In certain cases (Yuan et al., 2024; Shidani et al., 2025), synthetic data can even yield better results. Because it may better fit the model’s distribution and even provide regularization.	
1982		
1983		
1984		
1985		
1986		
1987		
1988		
1989		
1990		
1991		
1992		
1993	For these reasons, it is also a better choice for inductive reasoning.	
1994		
1995	<b>A.5 More Details about Benchmarks</b>	
1996	In this section, we will introduce the inputs and outputs of LLM inductive reasoning benchmarks one by one in detail.	
1997		
1998		
1999	<b>ILP</b> (Evans and Grefenstette, 2018). It takes as input background knowledge in first-order logic, along with a pair of positive and negative examples for a specific first-order logic case, and the model is required to output a first-order logic that satisfies these conditions.	
2000		
2001		
2002		
2003		
2004		
2005	<b>SCAN</b> (Lake and Baroni, 2018). It takes as input a series of entities and their states, and requires LLMs to output the actions needed to achieve those states. For example, given ‘a ball on the table’ as input, the model should output ‘place it on the table’.	
2006		
2007		
2008		
2009		
2010		
2011	<b>ARC</b> (Chollet, 2019). It takes as input several pairs of grids in natural language form, where they illustrate a specific transformation pattern. Then, given a new grid as input, the model is asked to output what the transformed grid would look like.	
2012		
2013		
2014		
2015		
2016	<b>List Functions</b> (Rule, 2020). It takes as input several pairs of number lists, where they illustrate a specific transformation pattern. Then, given a new number list as input, the model is asked to output what the transformed number list would look like.	
2017		
2018		
2019		
2020		
2021	<b>PROGES</b> (Alet et al., 2021). It provides several input-output pairs of a program and requires the model to generate the program itself.	
2022		
2023		
2024	<b>SyGuS</b> (Odena et al., 2021). It takes as input several pairs of strings, where they illustrate a specific transformation pattern. Then, it requires the model to generate a program to show this transformation process.	
2025		
2026		
2027		
2028		
	<b>ACRE</b> (Zhang et al., 2021a). It takes as input the results of interactions between different entities and a certain machine (i.e., the functions of different entities), and models need to output the entity corresponding to a specific function.	2029
		2030
		2031
		2032
		2033
	<b>Instructions</b> (Honovich et al., 2022). The model is given two natural language statements, A and B, where B is obtained by applying a certain instruction to A, and it is required to output that natural language instruction.	2034
		2035
		2036
		2037
		2038
	<b>Arithmetics</b> (Wu et al., 2024). It takes as input several pairs of two-digit additions along with their results, where these additions follow a calculation process in a certain numeral base. Given a pair of two-digit numbers, the model is required to output the result of their addition in the same base.	2039
		2040
		2041
		2042
		2043
		2044
	<b>Levy/Holt</b> (Liu et al., 2024b). It takes as input a pair of triplets, where the positive triplet represents a factual relationship and the negative triplet represents a counterfactual relationship. The task is to output an inference rule between triplets such that the positive triplet entails the negative triplet.	2045
		2046
		2047
		2048
		2049
		2050
	<b>NutFrame</b> (Guo et al., 2024). It takes text fragments that contain potential frames and frame elements, along with contextual information such as sentence structure, lexical cues, and semantic hints. This input helps the model identify underlying conceptual structures in the text. The model produces three types of outputs. Frame Induction: Identifies latent frames expressed in the text and maps them to existing FrameNet frames or proposes new frames. Frame Element Identification: Detects specific frame elements within the text. Frame Filling: Assigns concrete values from the text (entities or phrases) to the identified frame elements.	2051
		2052
		2053
		2054
		2055
		2056
		2057
		2058
		2059
		2060
		2061
		2062
		2063
	<b>DEER</b> (Yang et al., 2024). It takes as input several pairs of facts, where they illustrate a specific real-world rule. Then, it requires the model to generate the rule.	2064
		2065
		2066
		2067
	<b>RULEARN</b> (He et al., 2025) The input to it consists of three parts. Puzzle Scenarios: Each puzzle presents a set of conditions or operations that the agent can manipulate. These scenarios are designed to have underlying rules that are not explicitly provided. Agent Actions: The agent can perform a variety of actions, such as inputting integers or letters, to interact with the puzzle environment. Feedback: After each action, the agent receives feedback that	2068
		2069
		2070
		2071
		2072
		2073
		2074
		2075
		2076

2077	helps in refining its understanding of the hidden	<b>multi-step or compositional reasoning.</b> Many in-	2126
2078	rule. The desired output is: the hidden rule gov-	ductive tasks require composing several primitive	2127
2079	erning the puzzle scenario based on the feedback	operations (e.g., nested arithmetic, layered string	2128
2080	received from its actions.	transformations, or multi-step grid manipulations).	2129
2081	<b>Cryptography</b> (Li et al., 2025a) It takes as input	LLMs often fail to correctly chain these operations,	2130
2082	several pairs of English words, where each pair	producing partial reasoning steps that appear cor-	2131
2083	follows a certain cryptographic transformation pat-	rect individually but do not combine into a coherent	2132
2084	tern. Given a new word, the model is required to	global rule. (4) <b>Syntactic formatting errors.</b> The	2133
2085	output the new word obtained by applying the same	rule induced by the model may be semantically	2134
2086	transformation pattern.	correct, but it does not conform to the syntactic	2135
2087	<b>GeoILP</b> (Chen et al., 2025c) The same as ILP in	constraints of the task or the underlying logical	2136
2088	general.	system. For example, the generated code-like rule	2137
2089	<b>InductionBench</b> (Hua et al., 2025) It takes as	may contain syntax errors, making it impossible to	2138
2090	input a pair of strings, where the pair follows a	execute. (5) <b>Data sparsity challenge.</b> The rules	2139
2091	certain transformation rule, and the task is to output	induced by the model often apply only to the most	2140
2092	that rule.	frequent observations. When the distribution of	2141
2093	<b>CodeSeq</b> (Chen et al., 2025a) It takes as input	inductive patterns is uneven, the model is prone to	2142
2094	a number sequence of numbers and is required	errors on rare or boundary-case observations that	2143
2095	to output the number sequence’s general formula,	lie in the long tail.	2144
2096	with the entire output presented in code form.		
2097	<b>Statements</b> We determine the number of test	<b>A.7 Advantages of Sandbox and OC</b>	2145
2098	samples for each benchmark by consulting the orig-	<b>A.7.1 Advantages of Sandbox</b>	2146
2099	inal paper, the corresponding GitHub repository,	We will elaborate on the advantages of sandbox	2147
2100	and the Hugging Face dataset files, either through	testing from the following points.	2148
2101	direct counting or estimation.	(1) <b>Meeting the demand for complex real-</b>	2149
2102	<b>A.6 Failure Analysis</b>	<b>world tasks:</b> With the development of LLMs, there	2150
2103	Different benchmarks evaluate different models	is an increasing demand for handling more com-	2151
2104	and do not provide a standardized, unified evalu-	plex tasks in real-world scenarios, such as scientific	2152
2105	ation strategy. Therefore, for the sake of fairness,	discovery and simulations. Consequently, Agent	2153
2106	we cannot provide directly comparative results.	technologies have also advanced rapidly. In these	2154
2107	Hence, we summarize several typical failure	scenarios, task success rates cannot be adequately	2155
2108	modes of LLMs on these inductive tasks based	measured using traditional metrics. Sandbox test-	2156
2109	on prior experiences and the conclusions drawn in	ing, as a highly integrated evaluation framework	2157
2110	previous studies.	combining LLMs, tool usage, and complex rules, is	2158
2111	(1) <b>Difficulty in internalizing inductive logic.</b>	better suited to address modern AI tasks, including	2159
2112	LLMs tend to capture only surface-level associ-	inductive-related tasks. (2) <b>Isolation and safety:</b>	2160
2113	ations among observations, while failing to truly	Sandboxes are typically isolated environments with	2161
2114	acquire the underlying logical or mathematical re-	higher security, which helps reduce bias and risk	2162
2115	lations. For example, when presented with new	while allowing rapid problem detection. In contrast,	2163
2116	samples or longer samples following the same pat-	traditional metrics provide only coarse numerical	2164
2117	tern, the model often cannot induce the correct	indicators. (3) <b>Scalability:</b> Sandboxes are scalable,	2165
2118	rule or produce the correct answer. (2) <b>Spurious</b>	and existing sandbox frameworks can cover thou-	2166
2119	<b>pattern matching.</b> LLMs often latch onto acci-	sands of LLM tasks (Li et al., 2025c). (4) <b>Many</b>	2167
2120	idental or superficial patterns that happen to fit the	<b>works adopt a sandbox as the evaluation tool:</b>	2168
2121	limited observations, leading to spurious rules that	Their experimental results demonstrate that a sand-	2169
2122	do not reflect the true underlying structure. These	box can reveal a greater variety of errors (Milev	2170
2123	rules typically fail on even slightly varied test cases,	et al., 2025), provide more logical explainability	2171
2124	revealing that the induced pattern was an artifact	(Sheffler, 2025), and cover more boundary behav-	2172
2125	rather than a genuine abstraction. (3) <b>Failure in</b>	iors (Alshahwan et al., 2024).	2173
		Many existing LLM foundation trainings, such	2174
		as the Intern series, as well as the Minimax-M1	2175

and Seed-thinking models, employ large amounts of all kinds of reasoning data and sandbox validation to enhance the models’ fundamental reasoning capabilities.

### A.7.2 Advantages of OC

The OC evaluation method has the following advantages.

(1) **Overcoming averaging effects:** Metrics like ACC only measure overall accuracy, reflecting the average performance, but they do not reveal whether the model works effectively across all observations. Measuring OC provides a more fine-grained evaluation signal. (2) **Improved explainability:** This observation-level evaluation can uncover a greater variety of error types, allowing analysis of specific errors. (3) **Empirical support from existing work:** Previous studies (Chen et al., 2025a) (it reformulates number sequence as algorithmic problems and uses the OC of code solutions as the training signal) show that using observation-level supervision signals can effectively enhance a model’s inductive reasoning abilities.

## A.8 Future Works

We will further expand and discuss the prospects and potential applications of inductive reasoning in this subsection.

### A.8.1 In AI domain

**Inductive reasoning tasks, benchmarks, and data** The inductive reasoning capability of LLMs can be examined from various other perspectives, such as format imitation (Cheng et al., 2024a), cross-domain induction (Chen et al., 2024a), multimodal inductive (Wang et al., 2024b), and so on. Even a single data point has the potential to be expanded into either SFT or RL data. This serves as a concrete truth demonstrating the scalability of inductive reasoning data. Therefore, constructing more inductive reasoning datasets and evaluation methods is key to improving LLMs’ inductive reasoning capabilities.

**Enhancing the inductive reasoning ability of LLMs** We summarize several typical failure modes of LLMs on inductive tasks in Appendix A.6. Addressing these issues can further enhance the model’s inductive reasoning ability. Using simple and effective training data and architectures to identify universal inductive biases is a reasonable option. In Section 5, we conclude some analyses of inductive reasoning theoretically

from current practical experience. These analyses are highly beneficial for improving the inductive reasoning capabilities of LLMs. The paper also presents three major categories of methods, and it is expected that more methods will be proposed in the future.

Considering that BPE is the longest-established inductive reasoning task with the richest body of research, we will list several algorithms related to solving the BPE task below.

**Example I: BPE** Inductive program synthesis (or programming-by-example, PBE) is fundamentally an inductive reasoning task and one of the earliest to originate: given a few procedural input-output pairs, the system must generalize to unseen inputs. Li and Ellis (2024) argues that PBE is a highly general form of few-shot inductive inference, and they show that while pre-trained LLMs perform poorly on classic PBE tasks (e.g., list and string manipulations), fine-tuning on in-distribution PBE data substantially improves performance. To address the inherent lack of search in LLMs, Verbruggen et al. (2025) introduce a within-prompt search method: they sample lines of code, execute them on the input examples, feed back the semantic outcomes, and use that to guide further generation, effectively letting the model perform a semantic-aware search. Building on this inductive foundation, Lee et al. (2025b) propose a transductive framework: at test time, they actively select test inputs, have the LLM predict their outputs, and eliminate candidate programs that disagree — this refines the inductive hypothesis space using feedback. In addition, recent work such as Khan et al. (2025) breaks down complex PBE tasks into subtasks, with the LLM guiding the synthesis of these subtasks and then recombining them, thus improving the model’s ability to solve difficult problems. At the same time, the CodeARC (Wei et al., 2025) benchmark has been introduced to evaluate LLMs’ reasoning ability in real-world inductive synthesis settings: the model can interactively query a black-box target function, propose candidate functions, and iteratively refine them via a differential testing oracle.

These advances collectively show that under LLM-driven paradigms, PBE is being reactivated and reshaped: modern methods integrate fine-tuning, search, semantic execution, and increasingly also active learning and feedback mechanisms to boost generalization and robustness.

2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321  
2322  
2323  
2324  
2325  
2326

## A.8.2 In Real-world Domain

**Driving scientific discovery and innovation** In the context of scientific research, inductive reasoning is a vital method for uncovering natural laws. Artificial intelligence, especially generative AI, is becoming a new tool for scientists to propose hypotheses and design inductive experiments (Reddy and Shojaee, 2025). By learning patterns from large-scale data, these models can suggest plausible explanations and explore vast hypothesis spaces that would be difficult for humans to enumerate manually. This can serve a wide range of scientific fields, including the technological applications mentioned in Section 2.2.2.

**Example II: AI4Science** AI4Science constitutes a form of inductive reasoning by LLMs: instead of merely fitting data, LLMs generate symbolic hypotheses (programs) that generalize from observed examples, thereby performing induction over scientific phenomena. For instance, Shojaee et al. (2025) proposes new equation ‘skeletons’ using an LLM informed by scientific priors, and refines them with evolutionary search to discover mathematically and physically meaningful relations. This work touches on debates about the nature of scientific reasoning. Barnum (2012) argues that science is not purely induction, which raises interesting questions about whether LLM-based equation discovery should be framed as induction, deduction, or something more agentic. More broadly, Merler et al. (2024) uses LLMs to iteratively propose functions and optimize coefficients, reinforcing that LLMs are now capable of hypothesis generation rather than mere curve-fitting. Extensions like Wang et al. (2025a) further enrich this paradigm by combining data-driven insights with a reflective feedback loop for refining symbolic candidates. Taken together, these developments suggest that LLM-driven scientific discovery can be viewed as one of the future works of inductive reasoning.

### Human-machine collaboration in education

As cognitively-enhanced AI develops, future educational systems may evolve into "human-in-the-loop" intelligent agents (Chu et al., 2025). Inductive reasoning can help AI systems better understand students’ ways of thinking and their learning processes. By inferring latent patterns from students’ interactions and responses, such systems can adapt instruction strategies to individual needs and learning trajectories. Furthermore, inductive reasoning enables AI to provide more in-

terpretable feedback and personalized guidance, fostering more effective and engaging human-AI collaboration in education.

### New perspectives from philosophy and cognitive science

From the perspective of cognitive science, inductive reasoning is a crucial tool for understanding human thought and the structure of cognition. Future research may further explore the similarities and differences between the inductive reasoning of AI and that of human thinking (Trepczynski, 2024). Such investigations could shed light on whether AI systems rely on mechanisms analogous to human abstraction, generalization, and concept formation, or whether they achieve inductive competence through fundamentally different processes.

### Ethical and social implications

As inductive reasoning capabilities become more powerful, AI systems’ decisions or recommendations may embody inductive generalizations that ‘seem reasonable but could be wrong’. Avoiding inductive bias will thus become a significant issue (Ueno et al., 2022). Such biases may arise from spurious correlations, incomplete data, or mismatches between training environments and real-world contexts, leading to systematic errors that are difficult to detect. Addressing this challenge will require not only improved model design and evaluation but also greater transparency, uncertainty awareness, and human oversight to ensure that inductive inferences remain reliable and trustworthy.

2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357