THE ROLE OF DEDUCTIVE AND INDUCTIVE REASON-ING IN LARGE LANGUAGE MODELS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

031

032 033 Paper under double-blind review

Abstract

Large Language Models (LLMs) have achieved substantial progress in artificial intelligence, particularly in reasoning tasks. However, their reliance on static prompt structures, coupled with limited dynamic reasoning capabilities, often constrains their adaptability to complex and evolving problem spaces. In this paper, we propose the Deductive and InDuctive(DID) method, which enhances LLM reasoning by dynamically integrating both deductive and inductive reasoning within the prompt construction process. Drawing inspiration from cognitive science, the DID approach mirrors human adaptive reasoning mechanisms, offering a flexible framework that allows the model to adjust its reasoning pathways based on task context and performance. We empirically validate the efficacy of DID on established datasets such as AIW and MR-GSM8K, as well as on our custom dataset, Holiday Puzzle, which presents tasks about different holiday date calculating challenges. By leveraging DID's hybrid prompt strategy, we demonstrate significant improvements in both solution accuracy and reasoning quality, achieved without imposing substantial computational overhead. Our findings suggest that DID provides a more robust and cognitively aligned framework for reasoning in LLMs, contributing to the development of advanced LLM-driven problem-solving strategies informed by cognitive science models.

1 INTRODUCTION

"The measure of intelligence is the ability to change." – Albert Einstein

034 Large Language Models (LLMs), such as GPT-4, have transformed natural language processing 035 by excelling in tasks such as language translation, summarization, and question-answering (Ope-036 nAI, 2023), particularly in reasoning tasks and few-shot learning. However, there is ongoing debate 037 regarding their problem-solving reliability. According to Zhou et al. (2024), scaling up and finetuning LLMs enhances their capabilities but also diminishes reliability, introducing unpredictable 038 errors even in simple tasks and reducing the effectiveness of human supervision. Conversely, Li et al. (2024) highlights that the application of the Chain of Thought (CoT) (Wei et al., 2022b) 040 methodology significantly improves the accuracy of LLMs in arithmetic and symbolic reasoning 041 tasks by enabling inherently serial computations, which pose challenges for low-depth transformers. 042 Furthermore, Bubeck et al. (2023) observes that LLMs demonstrate a high degree of accuracy and 043 consistency in multi-step reasoning tasks, particularly when employing techniques such as CoT and 044 self-consistency (Wang et al., 2022). Additionally, reinforcement learning from human feedback (RLHF) has been shown to enhance model performance, notably reducing the incidence of harmful 046 or inaccurate outputs (Ouyang et al., 2022; Christiano et al., 2017). These insights suggest that, de-047 spite concerns related to model scalability and the potential for errors introduced during fine-tuning 048 (Zhou et al., 2024), LLMs can exhibit considerable reliability in complex reasoning tasks when 049 guided by structured methodologies and reinforced with human feedback. Ensuring the robustness of LLM outputs remains a critical priority, necessitating further investigation into strategies aimed 050 at enhancing model resilience and dependability. 051

Despite the notable success of LLMs, they face several limitations when dealing with more complex
 and evolving tasks. In particular, their rigidity in reasoning and difficulty in generalizing across
 diverse problem types present significant challenges. A key limitation of current LLMs is their

054 reliance on static prompt structures and patterns learned during training, which restricts their adaptability in novel or evolving contexts. 057 These models often apply fixed strategies to 058 problem-solving, leading to challenges in tasks that require logical reasoning, such as calculating family relationships, performing numer-060 ical comparisons, or counting specific charac-061 ters in a word (Nezhurina et al., 2024). Al-062 though these tasks may seem straightforward, 063 LLMs tend to depend on pre-learned patterns 064 instead of dynamically adjusting their reason-065 ing processes, resulting in errors in more com-066 plex problem spaces (Marcus, 2020; Hendrycks 067 et al., 2020). This inflexibility contrasts with 068 human problem-solving, which is typically it-069 erative and adaptive (Sloman, 2009). Humans use inductive reasoning to derive general rules 070 from specific instances and then apply deduc-071 tive reasoning to novel situations, allowing for 072 dynamic strategy adjustments based on task 073 complexity. In contrast, current LLMs lack this 074 level of flexibility in reasoning, which limits 075



Figure 1: Comparison of reasoning approaches in LLMs including the IO, CoT, and DID framework, highlighting the progression from direct output generation to dynamic inductive and deductive reasoning for more adaptive problemsolving.

their ability to generalize and adapt to more sophisticated scenarios. This underscores the need for increased attention and research in this area.

Moreover, LLMs have difficulty generalizing reasoning across tasks that require dynamic adjust-078 ment or incremental problem-solving. While LLMs can achieve high accuracy on specific tasks, 079 their performance often degrades when confronted with problems that evolve or require multi-step reasoning. This issue is especially evident in tasks where the model must balance different rea-081 soning strategies or integrate information from multiple sources. Tasks such as stepwise numerical 082 reasoning, temporal reasoning, or complex multi-step inference highlight gaps in LLMs' ability to 083 maintain consistent reasoning across different stages of a task. These models tend to produce incon-084 sistent or contextually inappropriate answers when required to adjust their reasoning dynamically 085 as the problem unfolds. The static and inflexible nature of their reasoning pipeline limits general-086 ization and accuracy, particularly when compared to human problem-solving, which adapts to new information in real-time. Despite improvements with techniques like CoT (Wei et al., 2022b), Tree-087 of-Thought (ToT) (Yao et al., 2024), Temperature-Tree-of-Thought (T²oT) (Cai et al., 2024), and 880 Graph-of-Thought (GoT) prompting (Besta et al., 2024), current LLMs still struggle to adjust their reasoning dynamically, resulting in difficulties in addressing more fluid and complex tasks. 090

091 To address these challenges, we propose the De-In-Ductive (DID) method, a novel approach de-092 signed to enhance LLM reasoning by integrating both inductive and deductive reasoning processes within the prompt construction framework as the Figure 1 shows. Grounded in cognitive science 093 models of human reasoning, the DID method enables LLMs to adjust their reasoning pathways dynamically in response to the task context and its evolving complexity. In the DID method, inductive reasoning is first used to derive general rules from specific instances, followed by deductive reason-096 ing to apply these rules in solving particular problems. This hybrid reasoning process mirrors human cognitive strategies, allowing the model to adjust its reasoning dynamically based on real-time feed-098 back. By employing this dynamic prompt strategy, the DID method improves the adaptability and 099 flexibility of LLMs, enabling them to better handle complex, evolving problem spaces. 100

We validate the effectiveness of the DID method on established benchmarks such as AIW and MR-GSM8K (Wei et al., 2022a; Cobbe et al., 2021), as well as our custom dataset, Holiday Puzzle, which includes tasks about holiday date calculations. By leveraging DID's hybrid prompt strategy, we observe significant improvements in both solution accuracy and reasoning quality, achieved without imposing substantial computational overhead. These results demonstrate the efficacy of DID in addressing the limitations of current LLM approaches. This work provides the following key contributions:

- We introduce a De-In-Ductive (DID) methodology that integrates both inductive and deductive reasoning within LLM prompt construction. This dynamic approach addresses key limitations of static prompt structures, enhancing the model's reasoning flexibility and adaptability.
 - Through empirical evaluations, we demonstrate that the DID method significantly enhances the adaptability and efficiency of LLMs across a diverse set of complex tasks. Moreover, DID improves solution accuracy and reasoning quality without incurring substantial computational overhead.

117 2 RELATED WORKS

119 **Cognitive Science and Deductive-Inductive Reasoning** Deductive and inductive reasoning are 120 foundational concepts in cognitive science for understanding human thought processes. Deductive 121 reasoning, formalized by philosophers like Kant, involves applying general principles to specific 122 cases, ensuring conclusions logically follow from premises. Inductive reasoning, conversely, gen-123 eralizes from specific observations to form broader conclusions, as highlighted in *The Riddle of* 124 Induction (Goodman, 1983). Cognitive models view these modes of reasoning as complementary, 125 where inductive reasoning generates hypotheses, and deductive reasoning tests them (Wason, 1960). 126 The interplay between these two reasoning methods has been shown to enhance problem-solving accuracy, particularly in uncertain domains where balancing exploration and validation is critical 127 (Johnson-Laird, 1983; Kahneman & Tversky, 1974). Research on mental models and heuristics 128 highlights how this dynamic combination allows for more flexible reasoning, especially in tasks 129 characterized by complexity or ambiguity. Problems in cognitive science are often classified based 130 on their structure and uncertainty: well-structured problems (e.g., mathematical proofs) lend them-131 selves to deductive reasoning, whereas ill-structured or open-ended problems (e.g., scientific discov-132 ery) require inductive reasoning to form plausible hypotheses from incomplete data (Funke, 2013). 133 Cognitive insights have been increasingly integrated into neural networks (L Griffiths et al., 2008), 134 with recent studies emphasizing the importance of embedding inductive structures within models to 135 improve generalization across tasks (Tenenbaum et al., 2011). The DID framework builds on this 136 cognitive science foundation by dynamically combining inductive and deductive reasoning within 137 LLMs, creating a hybrid model that mirrors human cognitive processes and enhances adaptability in problem-solving. 138

139

108

109

110

111

112

113

114

115

116

LLMs for Reasoning and Prompting Techniques While LLMs like GPT-4 have shown remark-140 able capabilities in tasks such as text generation and summarization (Brown et al., 2020), they often 141 struggle with structured reasoning, particularly in tasks involving logical inference, numerical com-142 parison, and complex deduction (OpenAI, 2023; Nezhurina et al., 2024). These shortcomings have 143 been well-documented in challenges like the ARC Prize (Rae et al., 2021) and in tasks that require 144 step-by-step reasoning or multi-hop inferences. To address these limitations, prompting techniques 145 such as CoT prompting (Wei et al., 2022b), ToT (Yao et al., 2024), and GoT (Besta et al., 2024) 146 have been developed to improve LLMs' capacity for structured reasoning. CoT enhances stepwise 147 reasoning by breaking down complex problems, while ToT and GoT explore multiple solution paths 148 through structured thought representations. However, these approaches remain static, requiring ex-149 tensive prompt engineering and lacking the dynamic adaptability needed for diverse tasks. Recent work, such as Hypergraph of Thoughts (HoT) (Yao et al., 2023), extends these methods to mul-150 timodal and more complex reasoning but still fails to offer the real-time adaptability needed for 151 nuanced problem-solving. Another critical analysis by Marcus (2020) highlights the limitations of 152 current LLMs, emphasizing their struggles with maintaining logical consistency and coherence in 153 complex reasoning tasks. The DID framework addresses these gaps by dynamically integrating in-154 ductive and deductive reasoning, making the reasoning process more adaptive and context-sensitive, 155 informed by the probabilistic reasoning advancements seen in the combination of Bayesian models 156 with neural networks (Gershman et al., 2015).

157 158

3 Methodology

- 159 160
- 161 The De-In-Ductive (DID) framework dynamically integrates inductive and deductive reasoning, inspired by cognitive science models of human reasoning. It balances inductive hypothesis generation

and deductive rule application to improve the flexibility and adaptability of large language models
 (LLMs) in complex problem-solving tasks.

3.1 PRELIMINARY

In this section, we formalize the core assumptions underlying the DID framework, which define the problem space and establish the performance limitations of LLMs in solving these problems. These assumptions lay the theoretical groundwork for modeling the interaction between problem complexity and the reasoning capabilities of LLMs.

Assumption 1 (Problem Distribution and Complexity) Let \mathcal{P} denote a problem space, where each problem instance $p \in \mathcal{P}$ is associated with an observed dataset $D_p = \{d_1, d_2, \dots, d_n\}$ and additional new data $D_{\text{new},p} = \{d'_1, d'_2, \dots, d'_m\}$. Each problem p is governed by a true underlying hypothesis H_p , which explains the relationship between the data and the problem context.

The observed data D_p follows a joint distribution conditioned on the hypothesis H_p :

$$P(D_p \mid H_p) = \prod_{i=1}^{n} P(d_i \mid H_p),$$
(1)

where d_i is independently generated according to the true hypothesis H_p . During the deductive phase, the model tests a candidate hypothesis H against the new data $D_{\text{new},p}$, and the likelihood is given by:

$$P(D_{\text{new},p} \mid H) = \prod_{j=1}^{m} P(d'_{j} \mid H).$$
(2)

The complexity of each problem p is characterized by a parameter $c(p) \in \mathbb{R}_+$, where higher values of c(p) indicate more complex problems, influencing the difficulty of hypothesis testing and data modeling.

Assumption 2 (Baseline Performance and Deductive Probability) Let \mathcal{M}_{LLM} represent a baseline LLM with parameters θ . For problems with complexity $c(p) \ge c_0$, the likelihood that \mathcal{M}_{LLM} produces a correct solution is bounded by:

$$\mathbb{P}_{\text{correct}}(p,\theta) \le \epsilon, \quad \text{for } c(p) \ge c_0, \tag{3}$$

where $\epsilon \ll 1$ reflects the baseline model's limitations on novel or complex tasks.

During the deductive reasoning phase, the likelihood that the hypothesis H holds, given the new data D_{new} , is computed by:

$$P_{\text{deductive}}(H \mid D_{\text{new}}) = \prod_{i=1}^{m} P(d'_i \mid H, \theta),$$
(4)

where the model tests the validity of H using the new observations D_{new} and adjusts its confidence in the hypothesis.

3.2 DE-IN-DUCTIVE (DID) FRAMEWORK

Figure 2 illustrates the comparison between the IO, CoT, and DID frameworks. The IO (Input-Output) Method processes natural language queries by retrieving patterns and facts without engaging
in iterative reasoning. The Chain of Thought (CoT) Method improves logical reasoning by breaking
down complex problems into sequential steps. Our proposed De-In-Ductive (DID) Method goes
further by dynamically integrating inductive and deductive reasoning. By iteratively generating and
testing hypotheses, DID adapts to problem complexities more effectively than static methods like



Figure 2: Comparison of reasoning approaches in Large Language Models (LLMs) including the IO
 method, Chain of Thought (CoT) prompting, and the De-In-Ductive (DID) framework, highlighting
 the progression from direct output generation to dynamic inductive and deductive reasoning for more
 adaptive problem-solving.

CoT, optimizing problem-solving by balancing reasoning modes in response to task difficulty. Based on Assumptions 3.1 and 3.1, DID optimizes an objective function that balances both reasoning modes, dynamically adjusting as the complexity of the problem increases.

Inductive Reasoning The inductive phase generates hypotheses from observed data D_p . Using Bayesian inference, the posterior probability for a hypothesis H is calculated as:

$$P(H \mid D_p) = \frac{P(D_p \mid H)P(H)}{P(D_p)},$$
(5)

where P(H) is the prior probability and $P(D_p | H)$ is the likelihood of the observed data under hypothesis H.

Deductive Reasoning In the deductive phase, the generated hypothesis H is tested against the new data $D_{\text{new},p}$. The likelihood of the new data given the hypothesis is:

$$P(D_{\text{new}} \mid H) = \prod_{j=1}^{m} P(d'_j \mid H).$$
(6)

This phase refines the hypothesis by comparing how well it explains the new observations.

Hybrid Objective Function The DID framework minimizes a hybrid objective function that integrates both inductive and deductive reasoning losses. The total objective function is defined as:

$$\mathcal{L}_{\text{DID}}(\theta) = \alpha \cdot \mathcal{L}_{\text{inductive}}(\theta) + (1 - \alpha) \cdot \mathcal{L}_{\text{deductive}}(\theta), \tag{7}$$

where $\mathcal{L}_{inductive}(\theta)$ represents the error in generalizing from the data, and $\mathcal{L}_{deductive}(\theta)$ penalizes errors in applying the rules to new instances. The weighting factor α adjusts dynamically based on the task complexity and the relative uncertainties in the inductive and deductive phases.

Dynamic Adjustment The balance between inductive and deductive reasoning is dynamically weighted based on the uncertainties in each reasoning process:

$$\alpha_t = \frac{\mathbb{E}_H[P_{\text{inductive}}(H \mid D_p)]}{\mathbb{E}_H[P_{\text{inductive}}(H \mid D_p)] + \mathbb{E}_{D_{\text{new}}}[P_{\text{deductive}}(D_{\text{new}} \mid H)]}.$$
(8)

This adaptive mechanism ensures that DID can adjust its strategy according to the evolving complexity of the problem.

277 278

273 274

275 276

Limitations of Chain-of-Thought (CoT) under Assumptions Under Assumptions 1 and 2, we infer that static reasoning methods such as Chain-of-Thought (CoT) prompting are insufficient for solving problems from the novel distribution \mathcal{P} . Since CoT relies on pre-learned patterns and fixed reasoning pathways, it lacks the dynamic adaptability required to handle new and complex problem structures. Specifically, without the ability to adjust α dynamically and integrate inductive hypothesis generation, CoT cannot effectively improve the low baseline performance ($\mathbb{P}(\text{Correct} \mid p, \theta) \leq \epsilon$) for complex problems.

Computational Complexity The computational complexity of the DID algorithm can be expressed as $O(T \cdot (n_{ind} + n_{ded}))$, where *T* is the number of iterations, n_{ind} is the complexity of inductive reasoning per iteration, and n_{ded} is the complexity of deductive reasoning per iteration. The dynamic adjustment enabled by the DID framework, as necessitated by Assumptions 1 and 2, reduces error propagation compared to static methods like CoT, thus improving efficiency and accuracy.

Efficiency Gains By dynamically adjusting reasoning pathways according to task complexity and
 problem novelty, DID reduces the number of iterations required to converge on a solution. This
 adaptive process ensures that DID outperforms static methods like CoT in both performance and
 computational overhead, achieving higher accuracy with fewer computational resources.

Theoretical adaption analysis The DID framework effectively manages cognitive load by initially focusing on simplified problem versions, allowing the model to concentrate on essential elements before engaging with more complex interactions. Through the inductive phase, the model observes specific instances, forming a foundation for generalization. As task complexity increases, the model transitions to deductive reasoning, applying generalized rules to arrive at a solution. This dynamic adjustment of reasoning strategies based on evolving task contexts enhances the model's adaptability and problem-solving efficiency.

305 Integration with Existing Models The De-In-Ductive (DID) method is compatible with various 306 LLM architectures and can be seamlessly integrated with existing techniques such as CoT prompting 307 (Wei et al., 2022b). By providing a structured reasoning framework that dynamically incorporates 308 both inductive and deductive reasoning, DID complements these methods and enhances the model's 309 problem-solving capabilities. Specifically, it structures the problem-solving process to allow for 310 the dynamic integration of reasoning strategies, utilizes a structured template prompt that guides 311 the model through incremental reasoning stages and improves adaptability and problem-solving 312 efficiency without introducing significant computational overhead. By mirroring human adaptive 313 reasoning processes, the DID method provides a more flexible and robust framework for LLMs to tackle complex and evolving problems. 314

315 316

317

4 EXPERIMENTS

318 4.1 ALICE PROBLEMS 319

Task. The AIW dataset is focused on evaluating the logical reasoning and deduction abilities of
 large language models (LLMs). The problems are structured around scenarios where models must
 infer relationships between family members based on a set of constraints, typically involving sib lings, with questions like determining how many sisters or brothers a particular sibling has. Due to
 the AIW GitHub open-source dataset being provided in the form of questions and various prompts,

4		Alice Problem			MR-GSM8K		Holiday Puzzle		
5	Model Prompt Method	IO (%)	CoT (%)	DID (%)	CoT (%)	DID (%)	IO (%)	CoT (%)	DID (%)
6	GPT-3.5 Turbo	6.7	8.6	13.3	68.1	73.3	0.2	1.4	5.6
20	GPT-40	43.4	55.9	70.3	82.0	83.7	7.8	5.2	15.4
27	Claude 3.5 Sonnet	74.8	83.7	89.5	91.3	92.0	17.4	17.8	24.5
28	·								

Table 1: Merged Results for GPT-3.5 Turbo, GPT-40, and Claude 3.5 Sonnet across Different Tasks (Alice Problem, MR-GSM8K, Holiday Puzzle)

we manually removed all prompts and eliminated duplicate questions that remained after prompt removal. This resulted in 113 unique original Alice problems. All subsequent experiments are based on these 113 problems, with results averaged over 20 runs.

Baseline and Framework Setup. We compare the performance of the DID framework with three other prompting methods: the IO prompt (which directly utilizes the LLM without structured prompting) and the CoT prompt. All methods are evaluated in a zero-shot setting. The comparisons are performed using three representative models: GPT-40, GPT-3.5-turbo, and Claude 3.5 Sonnet. GPT-40 was selected due to its highest reported accuracy in the AIW paper, while Claude 3.5 Sonnet achieved the second-highest accuracy. GPT-3.5-turbo, by contrast, demonstrated mid-to-low accuracy. For a fair comparison, all model parameters, including temperature, top-k sampling, and other hyperparameters, are maintained at their default values.



366

367

368

369

370

330

331 332

333

334

335 336

337

338

339

340

341

342

343

Figure 3: Comparison of reasoning approaches in Large Language Models (LLMs) including the IO method, Chain of Thought (CoT) prompting, and the De-In-Ductive (DID) framework, highlighting the progression from direct output generation to dynamic inductive and deductive reasoning for more adaptive problem-solving.

371 **Results** The results in 4 present a comparative analysis of the IO, CoT, and DID prompting meth-372 ods across three representative models: GPT-3.5 Turbo, GPT-4o, and Claude 3.5 Sonnet. Across all 373 models, the DID framework consistently outperforms both IO and CoT methods, demonstrating its 374 superior capability in handling multi-step reasoning tasks. 375

For GPT-3.5 Turbo, the DID method achieves an accuracy of 13.27%, significantly outperform-376 ing the IO prompt (6.73%) and CoT prompt (8.62%). This result highlights the DID framework's 377 ability to guide the model through structured, step-by-step reasoning, even in less powerful mod378 els like GPT-3.5 Turbo. Traditional prompt methods, such as IO and CoT, tend to struggle with 379 multi-step reasoning because they either rely on direct output generation or follow static, predefined 380 thought chains that lack the flexibility to adapt to more complex scenarios. On GPT-40, the DID 381 method reaches 70.27%, which far exceeds both IO (43.36%) and CoT (55.85%). The substantial 382 margin between DID and other methods demonstrates its efficacy in managing logical deductions, particularly in more complex reasoning tasks. In contrast, the static nature of IO and the limited 383 flexibility of CoT prevent these methods from fully addressing the intricacies of multi-step deduc-384 tions. Claude 3.5 Sonnet exhibits the highest overall performance across all methods, with the DID 385 method achieving an impressive accuracy of 89.49%, further extending its lead over IO (74.77%) 386 and CoT (83.68%). 387

388 As illustrated in Figure 3, the DID framework excels in progressively guiding LLMs through increasingly complex reasoning steps, especially in tasks where relationships must be deduced from 389 ambiguous information, such as determining family connections. Traditional prompt methods often 390 fail to handle such tasks effectively because they attempt to solve the original complex problem 391 directly, which can result in the model missing critical logical connections. This frequently leads 392 to incorrect conclusions. By breaking down multi-step reasoning problems into simpler subprob-393 lems, the DID method ensures that the model remains on track and avoids common pitfalls. This 394 structured approach enables LLMs to generalize effectively and handle more sophisticated logical 395 problems, like those presented in the Alice problem.

396 397 398

4.2 MR-GSM8K MATH PROBLEMS

Task. MR-GSM8K builds upon the GSM8K benchmark but introduces significantly higher complexity by focusing on meta-reasoning. Since MR-GSM8K not only challenges models to identify and explain errors in provided solutions, but also making it more suitable for assessing the advanced cognitive abilities of LLMs. Only the dataset portion of MR-GSM8K is used, in order to test different methods. And the assessing portion is not used. The dataset includes harder problem types, such as reversed reasoning and programmatic thinking, requiring deeper understanding and reasoning capabilities, thus offering a more rigorous evaluation framework for state-of-the-art models.

406

Baseline and Framework Setup. We compare the performance of the DID framework with the CoT prompt. The CoT is widely used as a method for improving the performance of LLM. All model parameters, including temperature, top-k sampling, and other hyperparameters, are set to their default values for a fair comparison.

Results Based on the results shown in 4, the performance comparison across different models (GPT-3.5 Turbo, GPT-4o, and Claude 3.5 Sonnet) for the CoT and DID frameworks reveals the consistent superiority of the DID framework. On GPT-3.5 Turbo, the DID method achieves an accuracy of 73.3%, outperforming CoT (68.1%). Similarly, on GPT-4o, the DID method demonstrates a clear advantage with an accuracy of 83.7%, compared to 82.0% for CoT. For Claude 3.5 Sonnet, the DID method further solidifies its dominance, achieving an accuracy of 92.0%, surpassing CoT (91.3%). This consistent performance across models highlights the effectiveness of the DID approach.

In summary, across GPT-3.5 Turbo, GPT-4o, and Claude 3.5 Sonnet, the DID framework consistently outperforms the CoT framework in terms of accuracy, demonstrating a significant performance advantage. This consistent superiority suggests that the DID method possesses stronger reasoning capabilities and higher reliability for handling complex tasks.

422 423

424

4.3 HOLIDAY PUZZLE

Task. The *Holiday Puzzle* is a task designed to evaluate the reasoning ability of LLMs in handling
holiday arrangements. This task includes 20 examples, each based on a holiday and compensatory
workday records from the past decade. The task requires the LLMs to calculate the actual number
of extra holiday days people receive, excluding weekends that were already scheduled as rest days.

429

Baseline and Framework Setup. To assess the effectiveness of the DID framework on the holiday
 calculation problem, we compare it with two alternative prompting methods: the IO prompt (which directly utilizes the LLM without structured prompting) and the CoT prompt. For fairness, all model



Figure 4: Comparison of reasoning approaches in Large Language Models (LLMs) including the IO
method, Chain of Thought (CoT) prompting, and the De-In-Ductive (DID) framework, highlighting
the progression from direct output generation to dynamic inductive and deductive reasoning for more
adaptive problem-solving.

parameters, including temperature, top-k sampling, and other hyperparameters, are kept at their default values.

463 **Results** As shown in 4, the performance comparison across different prompting methods (IO, CoT, 464 and DID) on the "Holiday Puzzle" task highlights the consistent superiority of the DID framework. 465 On GPT-3.5 Turbo, the DID method achieves an accuracy of 5.6%, significantly outperforming IO 466 (0.2%) and CoT (1.4%). This demonstrates that even on less powerful models, the DID framework effectively helps capture the underlying structure of complex temporal and scheduling relationships. 467 Similarly, on GPT-40, the DID method shows a clear advantage, with an accuracy of 15.4%, com-468 pared to 7.8% for IO and 5.2% for CoT. For Claude 3.5 Sonnet, the DID method solidifies its 469 dominance with an accuracy of 24.5%, far surpassing IO (17.4%) and CoT (17.8%). 470

471 In this task, the DID framework's stepwise combination of inductive and deductive reasoning proves 472 especially effective for capturing complex patterns such as the relationship *Holiday rest days* = *Total* 473 rest days - Weekend rest days. Figure 4 demonstrates how DID enables the model to approach this problem systematically, avoiding the risk of falling into incorrect reasoning paths common 474 in traditional prompt methods. By starting with simpler, deductively generated subproblems, the 475 DID framework ensures that the model can handle both the intricacies of date calculations and 476 generalize across varying input conditions. This step-by-step approach allows LLMs to refine their 477 reasoning process as complexity increases, ultimately improving their overall accuracy on tasks like 478 the Holiday Puzzle.

479 480 481

482

484

454

459 460

461

462

5 DISCUSSION

483 5.1 CHALLENGES OF ACHIEVING 100% ACCURACY IN SIMPLE TASKS

485 Despite the implementation of advanced prompting techniques such as CoT, ToT, and our proposed DID method, LLMs continue to face challenges in consistently achieving 100% accuracy, even on

486 seemingly simple logical tasks. A key factor underlying this limitation may be the fundamental 487 architecture of LLMs. These models rely on predicting the next token in a sequence, which restricts 488 their ability to maintain a coherent internal representation across multiple reasoning steps. While 489 attention mechanisms allow models to reference previous tokens, they lack the robust cognitive 490 structures that humans use to ensure logical integrity across reasoning processes. Consequently, LLMs can lose track of intermediate steps or overlook crucial logical connections, leading to errors 491 in tasks that might otherwise appear straightforward. This token-based, output-driven mechanism, 492 although effective in many natural language processing tasks, is inherently unsuited for tasks that 493 require rigorous logical consistency and structured reasoning, explaining the persistence of basic 494 mistakes in LLM outputs. 495

496 497

498

5.2 FINE-TUNING LLMs WITH DEDUCTIVE AND INDUCTIVE REASONING

499 The DID method proposed in this paper primarily focuses on prompting strategies without altering 500 the underlying architecture or fine-tuning the LLM itself. However, future work could investigate the benefits of fine-tuning LLMs on datasets explicitly incorporating deductive and inductive rea-501 soning processes. Fine-tuning strategies such as Reinforcement Learning from Human Feedback 502 (RLHF), Retrieval-Augmented Generation (RAG), or other methods (Lewis et al., 2020; Ouyang 503 et al., 2022; Christiano et al., 2017) could enhance the model's ability to handle complex reason-504 ing tasks. By integrating examples that demonstrate dynamic reasoning adjustments—similar to the 505 DID approach—during the fine-tuning phase, models would gain a deeper understanding of induc-506 tive and deductive reasoning patterns. Additionally, emerging techniques like Test-Time Training 507 (TTT) (Sun et al., 2024) could be explored to further improve models' adaptability and reasoning 508 performance during evaluation. These advancements are likely to enhance the consistency and relia-509 bility of LLM outputs in structured reasoning tasks by fostering more robust internal representations 510 of logical thought processes.

511

512 513 5.3 THE ARC PRIZE CHALLENGE

514 The ARC (Abstraction and Reasoning Corpus) Prize represents a particularly challenging bench-515 mark designed to evaluate an AI's capacity for abstract reasoning and generalization. Unlike con-516 ventional AI tasks that target pattern recognition based on fixed datasets, ARC tests a model's ability 517 to generalize across diverse problem types, making it an ideal platform to assess the adaptability of 518 frameworks such as our De-In-Ductive (DID) method. ARC problems require abstract reasoning 519 across domains like visual pattern recognition and symbolic problem-solving-areas where DID's 520 integration of inductive and deductive reasoning could provide a competitive edge. The DID framework's ability to adapt dynamically based on task complexity, first using inductive reasoning to 521 hypothesize patterns and then applying deductive logic to test them, positions it well for tasks in the 522 ARC benchmark. Given that ARC problems often involve iterative problem-solving and the ability 523 to generalize from minimal examples, we believe the DID method is particularly suited for this task. 524 Future work will focus on evaluating DID's performance on ARC Prize tasks to test its robustness 525 and effectiveness in abstract reasoning scenarios. 526

520 527

6 CONCLUSION

528 529 530

In this work, we introduced the De-In-Ductive (DID) method, a novel framework that dynamically 531 integrates inductive and deductive reasoning to enhance the adaptability and reasoning capabilities 532 of Large Language Models (LLMs). By leveraging cognitive science principles, the DID framework 533 allows LLMs to evolve their problem-solving strategies in response to task complexity, overcoming 534 the rigidity of static prompt structures. Through extensive empirical validation on both standard benchmarks and our custom Holiday Puzzle dataset, we demonstrated significant improvements 536 in accuracy and reasoning quality, achieved without excessive computational costs. However, while 537 DID advances the field, challenges remain in making LLMs more intelligent, particularly in ensuring better generalization to unseen tasks, maintaining adaptability in complex multi-step reasoning, and 538 further refining model biases. Future research must continue to address these issues, paving the way for more robust and cognitively aligned artificial intelligence systems.

540 REFERENCES

549

562

568

542	Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gian-
543	inazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of
544	thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI
545	Conference on Artificial Intelligence, volume 38, pp. 17682–17690, 2024.

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- ⁵⁵³ Chengkun Cai, Xu Zhao, Yucheng Du, Haoliang Liu, and Lei Li. *t*² of thoughts: Temperature tree
 elicits reasoning in large language models. *arXiv preprint arXiv:2405.14075*, 2024.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lasse Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Joachim Funke. Complex problem solving: A case for complex cognition? In *Complex problem* solving: Principles and mechanisms, pp. 25–47. Psychology Press, 2013.
- Samuel J Gershman, Eric J Horvitz, and Joshua B Tenenbaum. Computational rationality: A converging paradigm for intelligence in brains, minds, and machines. *Science*, 349(6245):273–278, 2015.
- 569 Nelson Goodman. Fact, Fiction, and Forecast. Harvard University Press, Cambridge, MA, 1983.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Philip Nicholas Johnson-Laird. *Mental models: Towards a cognitive science of language, inference, and consciousness.* Number 6. Harvard University Press, 1983.
- Daniel Kahneman and Amos Tversky. Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131, 1974.
- 579 Thomas L Griffiths, Charles Kemp, and Joshua B Tenenbaum. Bayesian models of cognition. 2008.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.

- Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024.
- Gary Marcus. The next decade in ai: four steps towards robust artificial intelligence. *arXiv preprint arXiv:2002.06177*, 2020.
- Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. Alice in wonderland:
 Simple tasks showing complete reasoning breakdown in state-of-the-art large language models.
 arXiv preprint arXiv:2406.02061, 2024.

OpenAI. GPT-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

594 595 596 597	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35: 27730–27744, 2022.
598 599 600 601	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> , 2021.
602 603	S Sloman. <i>Causal models: How people think about the world and its alternatives</i> . Oxford University Press, 2009.
604 605 606 607	Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, et al. Learning to (learn at test time): Rnns with expressive hidden states. <i>arXiv preprint arXiv:2407.04620</i> , 2024.
608 609	Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. How to grow a mind: Statistics, structure, and abstraction. <i>science</i> , 331(6022):1279–1285, 2011.
610 611 612 613	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh- ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> , 2022.
614 615	Peter C. Wason. On the failure to eliminate hypotheses in a conceptual task. <i>Quarterly Journal of Experimental Psychology</i> , 12(3):129–140, 1960.
616 617 618	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Andy Yu, Brian Lester, Xuezhi Li, Yuan Cao, et al. Emergent abilities of large language models. <i>arXiv</i> preprint arXiv:2206.07682, 2022a.
619 620 621 622	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022b.
623 624 625	Fanglong Yao, Changyuan Tian, Jintao Liu, Zequn Zhang, Qing Liu, Li Jin, Shuchao Li, Xiaoyu Li, and Xian Sun. Thinking like an expert: Multimodal hypergraph-of-thought (hot) reasoning to boost foundation modals. <i>arXiv preprint arXiv:2308.06207</i> , 2023.
626 627 628	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
630 631 632	Lexin Zhou, Wout Schellaert, Fernando Martínez-Plumed, Yael Moros-Daval, Cèsar Ferri, and José Hernández-Orallo. Larger and more instructable language models become less reliable. <i>Nature</i> , pp. 1–8, 2024.
633 634	
635 636	
637 638	
639	
640	
641 642	
643	
644	
645	
646	
647	