Logical Reasoning in Large Language Models: A Survey

Anonymous ACL submission

Abstract

With the emergence of advanced reasoning models like OpenAI o4 and DeepSeek-R1, large language models (LLMs) have demonstrated remarkable reasoning capabilities. However, their ability to perform rigorous logical reasoning remains an open question. This survey synthesizes recent advancements in logical reasoning within LLMs, a critical area of AI research. It outlines the scope of logical reasoning in LLMs, its theoretical foundations, and the benchmarks used to evaluate reasoning proficiency. We analyze existing capabilities across different reasoning paradigms deductive, inductive, abductive, and analogical - and assess strategies to enhance reasoning performance, including data-centric tuning, reinforcement learning, decoding strategies, and neuro-symbolic approaches. The review concludes with future directions, emphasizing the need for further exploration to strengthen logical reasoning in AI systems.

1 Introduction

002

007

011

013

017

019

033

037

041

Logical reasoning is a fundamental challenge to artificial intelligence (AI) and natural language processing (NLP) (Newell and Simon, 1956; Mc-Carthy and Hayes, 1981; McCarthy, 1959). While early formal logic-based reasoning approaches faced limitations in scalability and adaptability (Pereira, 1982; Cann, 1993), data-driven models became the dominant method since the 1980s (Mc-Carthy, 1989). Recently, pre-trained Large Language Models (LLMs) and their emergent logical reasoning abilities have attracted increasing attention (Liu et al., 2023b; Xu et al., 2023). Logical reasoning integrates LLMs with inference structuring, enabling multistep deduction and abstraction, and improving interpretability and reliability (Shi et al., 2021; Stacey et al., 2022; Rajaraman et al., 2023). It also strengthens generalization, helping models handle novel scenarios beyond their

training data (Haruta et al., 2020). As LLMs become integral to domains like legal analysis and scientific discovery, ensuring the correctness and verifiability of their reasoning is increasingly vital. As a result, post-training LLM for reasoning has garnered a surge of interest in both industry and research(OpenAI, 2024; DeepSeek-AI, 2025; Muennighoff et al., 2025). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Recent surveys have touched upon LLMs' reasoning (Li et al., 2025; Huang and Chang, 2023). However, existing surveys discuss general reasoning, exemplified by chain-of-thought (CoT), treating logical reasoning as a task case, without dedicated discussion. There has been lack of a thorough literature review focusing on LLMs and formal symbolic logic. To address this issue, this survey provides a comprehensive review of logical reasoning in large language models (LLMs), with a focus on formal and symbolic logic-based reasoning rather than general heuristic approaches. The structure is illustrated in Figure 1. We begin by defining logical reasoning in AI, distinguishing it from general-purpose reasoning, and categorizing key paradigms, including deductive, inductive, abductive, and analogical reasoning. Then, we analyze existing benchmarks and evaluation methodologies, identifying gaps in assessing symbolic inference, consistency, and robustness. We further explore state-of-the-art techniques for enhancing logical reasoning, such as supervised fine-tuning, logic-informed pre-training, reinforcement learning, inference-time decoding strategies, and hybrid neuro-symbolic methods. We examine recent advances in neuro-symbolic integration, along with applications of theorem provers, logic solvers, and formal verification frameworks in LLMs. Finally, we highlight open challenges in scalability, reasoning consistency, explainability, and efficiency, proposing future directions for multi-modal reasoning, hybrid architectures, and improved evaluation frameworks.



Figure 1: The structure of this survey

2 Logic in Artificial Intelligence

083

087

091

100

101

Logical reasoning is a cornerstone of artificial intelligence (AI), enabling machines to simulate human thought processes and solve complex problems. At its core, logical reasoning applies structured rules to derive conclusions from premises, providing a rigorous framework for decision-making and inference (Sun et al., 2023).

2.1 History of Logic Reasoning Research

Logical reasoning can be traced back to ancient Greece, where Aristotle's syllogisms laid the foundation for classical logic. During the Middle Ages, scholars refined these theories, and in the 17th century, Leibniz's universal language and calculus ratiocinator bridged logic with mathematics, foreshadowing modern computational logic. The 19th century saw George Boole's Boolean algebra, which transformed logic into a mathematical framework, laying the foundation for digital computing.

The 20th century ushered in modern logic, with 102 Russell and Whitehead's Principia Mathematica 103 formalizing complex logical systems. By the midcentury, AI pioneers like John McCarthy leveraged 105 logic for knowledge representation and automated 106 theorem proving, leading to logic programming and knowledge bases. The 1970s introduced non-108 109 monotonic logic, enabling AI to handle commonsense reasoning. The 1980s saw logical reasoning 110 integrate with knowledge representation, advanc-111 ing expert systems for real-world applications. The 112 1990s saw the rise of knowledge graphs, structuring 113

vast knowledge for complex reasoning tasks.

With the development of deep learning in the 21st century, neuro-symbolic approaches stand out as a new approach for combining deep learning with logical inference, resulting in tools like Deep-Logic (Cingillioglu and Russo, 2019) and SAT-Net (Wang et al., 2019). Logical reasoning remains a cornerstone of AI research, evolving from philosophy to modern computing. As AI advances, logical reasoning continues to shape intelligent systems, ensuring structured, interpretable, and robust decision-making.

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

2.2 Types of Logical Reasoning

Logical reasoning can be broadly categorized into four main types, each serving distinct purposes and applications:

Deductive Reasoning. This type of reasoning derives specific conclusions from general principles or premises. It operates under the rule that if all premises are true and the reasoning is valid, the conclusion must also be true. Deductive reasoning is fundamental in fields such as mathematics and formal logic, where certainty and rigor are paramount.

Inductive Reasoning.Unlike deductive reason-
ing, inductive reasoning draws general conclusions138based on specific observations or evidence.140the conclusions are often considered probable, they
are not guaranteed to be true.141ure not guaranteed to be true.142widely used in scientific discovery and data-driven143

Dataset	Lang.	Task	Reasoning Type	Size	Source
LOGIQA	Zh/En	MRC	Misc.	8,678	Exam
ReClor	En	MRC	Misc.	6,138	Exam
AR-LSAT	En	MRC	Misc.	2,064	Exam
CLUTRR	En	MRC	Inductive	6,016	Rule
GSM	En	MRC	Deductive	19K	Exam
LINGOLY	En	MRC	Inductive	1,133	Expert
ConTRoL	En	NLI	Deductive	8,325	Exam
FOLIO	En	NLI	Deductive	1,351	Expert
LogicNLI	En	NLI	Deductive	30K	Exam
ProofWriter	En	NLI	Deductive	-	Exam
LogicBench	En	NLI	Deductive	1,270	Rule
ART	En	NLI	Abductive	20K	Expert
Analogical	En	NLI	Analogical	720K	Crowd
GLoRE	Zh/En	Misc.	Misc.	17 tasks	Misc.
LogiGLUE	En	Misc.	Misc.	24 tasks	Misc.
LogiTorch	En	Misc.	Misc.	16 tasks	Misc.
BIG-Bench	En	Misc.	Misc.	7 tasks	Misc.

Table 1: Main Datasets and Benchmarks of Logical Reasoning Task.

decision-making, where patterns and trends are inferred from empirical data.

Abductive Reasoning. This form of reasoning seeks the most plausible explanation or cause for a set of observations, often in the presence of incomplete information. Abductive reasoning is particularly useful in diagnostic tasks and real-world problem-solving. While abductive conclusions are not certain, they provide a practical basis for hypothesis generation and decision-making under uncertainty.

Analogical Reasoning. Analogical reasoning involves drawing comparisons between similar situations or domains to make inferences or solve problems. By identifying parallels between different scenarios, this type of reasoning enables creative problem-solving and knowledge transfer. Analogical reasoning is particularly valuable in fields like education, design, and innovation.

3 Tasks and Benchmarks

Logical reasoning datasets and benchmarks are essential for evaluating the reasoning capabilities of large language models (LLMs). These datasets can be categorized into three types based on their data sources:

Rule-based Datasets (Tafjord et al., 2021; Sinha et al., 2019) are automatically generated using logical rules, enabling large-scale data collection. However, ensuring diversity is crucial to avoid repetitive patterns and comprehensively evaluate reasoning capabilities.

Crowdsourced Datasets (Bhagavatula et al., 2020) leverage collective human intelligence for diverse reasoning tasks. While requiring quality

control measures, these datasets capture nuanced linguistic patterns and commonsense knowledge that automated methods often miss.

Expert-Designed Datasets (Han et al., 2024a) are constructed by domain experts, ensuring high precision and accuracy. Although typically smaller than crowd-sourced corpora, their meticulous design makes them indispensable for in-depth logical reasoning evaluation.

Exam-Based Datasets (Liu et al., 2021b; Yu et al., 2020; Wang et al., 2022) originate from standardized test questions (e.g., Chinese National Civil Service Exam, LSAT, CAT), offering high-quality, expert-crafted logic problems at scale. These datasets are widely used to evaluate reasoning in real-world scenarios.

Table 1 summarizes important datasets for logical reasoning, which typically cover tasks such as Natural Language Inference (NLI) (§3.1), Machine Reading Comprehension (MRC) (§3.2). Examples can be found in A.1.

3.1 Natural Language Inference (NLI)

NLI evaluates whether a *hypothesis* logically follows from a *premise*, directly assessing a model's reasoning ability. Labels typically fall into binary (Entailment, Non-entailment) or ternary (Entailment, Contradiction, Neutral) classifications. Some datasets use True and False labels instead. ConTRoL (Liu et al., 2021a), derived from recruitment exams, contains 8,325 entries labeled as Correct, Incorrect, or Can't Say, corresponding to entailment, contradiction, and neutral. FO-LIO (Han et al., 2024a), an expert-constructed dataset for First-Order Logic (FOL) reasoning, consists of 1,351 entries labeled True or False. LogicNLI (Tian et al., 2021), with 30K entries generated via logical rules, isolates FOLbased inference from commonsense reasoning. ProofWriter (Tafjord et al., 2021) extends Rule-Taker (Clark et al., 2021) by introducing closedworld (CWA) and open-world (OWA) assumptions, covering handcrafted domain theories and crowdsourced paraphrased rules for linguistic and domain generalization. LogicBench (Parmar et al., 2023), generated by GPT-3, includes 1,270 test entries across 25 reasoning types (e.g., propositional logic, FOL) labeled Yes or No. ART (Bhagavatula et al., 2020) contains 20K commonsense narrative contexts and 200k explanations for abductive reasoning evaluation.

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

178

179

180

181

145 146

144

- 147
- 148 149
- 150 151 152
- 153 154
- 155 156
- 157 158
- 159
- 160 161
- 162

165

166

167

168

169

170

171

172

173

174

175

176

230

231

235

240

241

243

244

245

246

247

248

254

258

263

265

267

272

273

274

276

3.2 Machine Reading Comprehension (MRC)

Machine Reading Comprehension (MRC) evaluates logical reasoning by requiring models to answer questions based on a given passage, commonly formatted as multiple-choice, span extraction, or free response. LogiQA (Liu et al., 2023a), sourced from the Chinese Civil Service Exam, contains 15,937 Chinese and English entries targeting complex logical reasoning. ReClor (Yu et al., 2020), derived from the GMAT, features 6,138 English multiple-choice questions with four options. AR-LSAT (Wang et al., 2022), collected from the LSAT exam, includes 2,064 entries covering ordering, grouping, and allocation games with five options each. CLUTRR (Sinha et al., 2019) focuses on inductive reasoning for kinship relationships in short narratives, containing 6,016 entries combining entity extraction and logical inference. LINGOLY (Bean et al., 2024) uses Linguistic Olympiad puzzles (1,133 problems across 6 formats and 5 difficulty levels) to assess pattern identification and generalization in low-resource or extinct languages.

3.3 Multi-modal Logical Reasoning

Recently, logical reasoning combining texts and images has been explored in several research. **LogicVista** (Xiao et al., 2024) collects 448 comprehensive exam-based logical reasoning data in Visual contexts with human-annotated rationales suitable for both open-ended and multiple-choice evaluation. **VISUALPUZZLES** (Song et al., 2025) is a holistic benchmark of 1,167 exam-based puzzlelike multi-modal questions specifically designed to decouple reasoning abilities from domain knowledge.

3.4 Benchmark Suites

Benchmark suites standardize evaluation and facilitate model comparison in logical reasoning research. GLoRE (liu et al., 2023) provides 17 test-only datasets for few-shot and zero-shot evaluation of generalization capabilities. LogiGLUE (Luo et al., 2024) unifies 24 logical reasoning tasks into a sequence-to-sequence format with both training and test sets for comprehensive evaluation. Logi-Torch (Helwe et al., 2022) offers a PyTorch-based framework with 16 datasets and model architectures for streamlined logical reasoning experiments.
BIG-bench (Srivastava et al., 2022) includes 7 collaborative logical reasoning tasks such as Logic

Grid Puzzle and Logical Fallacy Detection. LogiEval and LogiEval-Hard (Liu et al., 2025) provides a holistic testing suite with various sub-tasks for logical reasoning. 277

278

279

281

282

283

284

285

286

287

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

321

322

323

4 Evaluations

The rapid development of pre-trained language models (PLMs) necessitates rigorous evaluation of their logical reasoning capabilities. This section examines four reasoning paradigms—deductive, inductive, abductive, and analogical—while analyzing evaluation approaches and metrics.

4.1 Deductive Reasoning

Deductive reasoning, deriving specific conclusions from general premises, is crucial for automated theorem proving. Despite LLMs performing well on tasks like compositional proofs, standard benchmarks and encoding entailment relationships, they struggle with extended reasoning, hypothetical subproofs without examples, generalization, and sensitivity to syntactic variations (Saparov et al., 2023; Yuan et al., 2023; Ryb et al., 2022).

4.2 Inductive Reasoning

Inductive reasoning, which generalizes from specific instances to broader rules, is essential for tasks like hypothesis generation and pattern recognition. While Yang et al. (2024b) find that pre-trained models can serve as effective "reasoners", Bowen et al. (2024) show that even advanced LLMs struggle with simple inductive tasks in their symbolic settings. Similarly, Sullivan (2024) demonstrates that Transformer models, even after fine-tuning, fail to learn fundamental logical principles, indicating limited inductive reasoning capabilities.

4.3 Abductive Reasoning

Abductive reasoning, which seeks the most plausible explanations for observed phenomena, is crucial in fields like law and medicine. Del and Fishel (2023) highlights the challenges LLMs face in generating plausible hypotheses from incomplete information. In the legal domain, Nguyen et al. (2023) show that despite strong performance, models struggle with abductive reasoning, underscoring the complexity of this paradigm.

4.4 Analogical Reasoning

Analogical reasoning, which infers unknown information by comparing it with known information, is vital for tasks requiring creativity and knowledge transfer. Wijesiriwardene et al. (2023) introduced ANALOGICAL, a benchmark for long-text analogical reasoning. They find that as analogy complexity increases, LLMs struggle to recognize analogical pairs. Petersen and van der Plas (2023) show that models can learn analogical reasoning with minimal data, approaching human performance. However, Qin et al. (2024) question whether LLMs truly rely on analogical reasoning, discovering that random examples in prompts often achieve comparable performance to relevant examples.

4.5 Overall Analysis and Metrics

324

325

329

330

333

334

336

340

341

342

343

345

347

364

372

Liu et al. (2023b) evaluate GPT-4 and ChatGPT on benchmarks like LogiQA and ReClor, showing that while GPT-4 outperforms ChatGPT, both of them struggle with out-of-distribution tasks. Xu et al. (2023) introduce the NeuLR dataset and propose a framework evaluating LLMs across six dimensions: correctness, rigor, self-awareness, proactivity, guidance, and absence of hallucinations.

Metrics for Evaluating Logical Reasoning. Reasoning is fundamentally process-oriented rather than outcome-oriented (Leighton, 2003). Although traditional conclusion-based metrics like accuracy and F1 score are widely used for their simplic-348 ity and general applicability, they fall short in assessing the logical reasoning process (Mondorf and Plank, 2024). Recent studies have introduced rationale-based metrics to evaluate the reasoning trace. Structural parsing approaches (Saparov et al., 2023; Dziri et al., 2023) decompose the reasoning 354 process into formalized representations or graphs to facilitate a more fine-grained evaluation. However, these methods are typically constrained by 358 the requirement for specially structured reasoning texts. Other researchers have proposed interpretable quantitative metrics (Golovneva et al., 2022; Prasad et al., 2023), designing a variety of indicators to assess diverse properties of model reasoning. Nevertheless, these methods often rely on complex feature engineering and typically use metrics such as BERTScore or entropy, whose values lack clear physical interpretability. There remains a pressing need for a widely accepted and general rationale-based evaluation method.

5 **Enhancement Methods**

Enhancing LLMs' logical reasoning remains crucial. This section focuses on core strategies: Data-Centric Approaches (§5.1), Model-Centric Approaches (§5.2), External Knowledge Utilization 373 (§5.3), and Neuro-Symbolic Reasoning (§5.4). 374

5.1 **Data-Centric Approaches**

Data-centric approaches enhance LLMs' reasoning capabilities by utilizing meticulously curated training datasets. Formally:

$$D^* = \arg\max_D R(M_D) \tag{1}$$

375

376

377

378

383

384

385

386

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

where:

- D: training datasets.
- M_D : model trained on D.
- R: performance evaluator.

This formulation highlights the central role of dataset optimization in data-centric approaches. In practice, data-centric methods typically involve three types of datasets: expert-curated datasets, synthetic datasets, and LLM-distilled datasets.

Expert-Curated Datasets. The FOLIO series (Han et al., 2024a,b) establish formal verification through FOL annotations, with P-FOLIO extending the complexity of reasoning chains for enhanced training. LeanDojo (Yang et al., 2023) provides 98k+ human-proven mathematical theorems. Additionally, Symbol-LLM (Xu et al., 2024a) systematically organizes 34 symbolic reasoning tasks to capture inter-symbol relationships across 20 distinct symbolic families. These datasets benefit from high-quality data, with the scale commonly limited by labor-intensive annotation.

Rule-based Synthetic Datasets. Rule-based synthetic data remains fundamental for data generation. RuleTaker (Clark et al., 2021) formalizes this through a three-phase pipeline: behavior formalization, example synthesis and linguistic equivalents generation. Similarly, Morishita et al. (2024) develops Formal Logic Deduction Diverse (FLD $_{\times 2}$), a synthetic dataset based on symbolic theory and previous empirical insights. Rule-based data generation enables large-scale, systematic data creation and fine-grained control over specific reasoning patterns. However, rule-based datasets often suffer from limited linguistic diversity, reduced realism and a gap between artificial templates and realworld inference complexity.

LLM-Distilled Datasets. Researchers employ advanced models such as GPT-4 and Deepseek-R1 for intermediate reasoning step distillation.

505

506

507

508

509

510

511

512

513

514

515

465

LogiCoT (Liu et al., 2023c) augments existing 419 420 datasets with GPT4-generated reasoning chains, while LogicPro (Jiang et al., 2024) combines al-421 gorithmic problems with code solutions to create 422 variable-guided reasoning data. To advance, Wang 423 et al. (2024b) propose PODA, which generates 424 contrastive analyses of correct/incorrect options 425 through premise-oriented augmentation, enabling 426 reasoning path differentiation via contrastive learn-427 ing. These methods leverage the reasoning ability 428 of advanced language models and tackle the lack 429 of diversity and complexity of rule-based meth-430 ods. However, training models on LLM-distilled 431 data may lead to model collapse (Shumailov et al., 432 2023), causing the models to lose diversity and ac-433 curacy due to the loss of long-tail data present in 434 real distributions. 435

5.2 Model-Centric Approaches

Model-Centric approaches enhance LLMs' reasoning capabilities by optimizing model parameters and decoding strategies. The formal objective is:

$$(\theta^*, S^*) = \arg\max_{\theta, S} R(M_\theta, S)$$
(2)

where:

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

• M_{θ} : model with learnable parameters θ .

- S: decoding strategy (e.g., chain-of-thought prompting, verification-based decoding).
- R: reasoning performance metric.

This formulation highlights the joint optimization of model parameters θ and decoding strategy S. Practical implementations can be categorized as:

- Instruction Fine-Tuning: optimizing θ .
- Reinforcement Learning: optimizing θ .
- Inference-Time Decoding: optimizing S.

Model-Centric approaches focus on directly improving models' reasoning capabilities by optimizing its internal mechanisms and decoding strategies, making them complementary to data-centric approaches.

5.2.1 Supervised Fine-Tuning

Supervised Fine-Tuning (SFT) optimizes LLMs
through supervised learning on pairs of inputs and
desired outputs. For example, Liu et al. (2023c)
designs multi-grained logical instructions spanning
diverse levels of abstraction and complexity. Similarly, Feng et al. (2024) SFT models to mimic

logical solvers by replicating formal deduction reasoning processes. In addition, Xu et al. (2024a) implements two-stage symbolic fine-tuning through *Injection* (injecting symbolic knowledge) and *Infusion* (balancing symbol and NL reasoning).

To overcome SFT's over-fitting limitations, Wang et al. (2024b) enforce contrastive learning between factual/counterfactual paths with SFT. Further, Wang et al. (2024a) augments Llama models with a Program-Guided Learning Framework and logic-specific architecture adjustments.

In summary, the primary purpose of SFT in logical reasoning is generally to inject into the model the capability for specific reasoning manners such as symbolic reasoning or long CoT reasoning. Yue et al. (2025) demonstrates that SFT is able to introduce new reasoning patterns that are not present in the base model. However, high-quality SFT data relies on distillation from more advanced models or human annotation, which is more general and effective for expanding smaller models' boundaries. For advanced large models, obtaining and scaling up higher-quality SFT data is often challenging.

5.2.2 Reinforcement Learning

Reinforcement learning (RL) has become pivotal in optimizing large language models (LLMs), particularly since the breakthrough of Reinforcement Learning from Human Feedback (RLHF). Jiao et al. (2024) leverage RL for planning-based reasoning optimization, while Xi et al. (2024) develop R^3 , achieving process supervision benefits through outcome-only supervision.

The success of large-scale RL in OpenAIol (OpenAI, 2024) has inspired numerous studies. RL algorithms train ol-style models to enhance Chain-of-Thought (CoT) reasoning, addressing issues like formulaic outputs and limited long-form reasoning. For instance, Zhao et al. (2024) integrates CoT instruction fine-tuning with Monte Carlo Tree Search (MCTS) decoding for multipath reasoning exploration. In contrast, Zhang et al. (2024) employs MCTS to generate code-reasoning data for supervised fine-tuning (SFT) and Direct Preference Optimization (DPO).

A significant breakthrough comes from DeepSeek-R1 (DeepSeek-AI, 2025), which pioneers a novel RL strategy to enhance logical reasoning. DeepSeek-R1-Zero, trained purely through RL without SFT, demonstrates impressive reasoning capabilities but faces challenges in readability and language consistency. To address this,

DeepSeek-R1 introduces minimal long-CoT SFT 516 data as a cold start before RL, achieving a balance 517 between usability and reasoning performance. By 518 iteratively synthesizing high-quality reasoning data 519 through RL, DeepSeek-R1 overcomes limitations 520 imposed by human annotators, addressing issues 521 such as mechanistic responses, repetitive patterns, 522 and insufficient long-chain reasoning. This approach represents a potential paradigm shift in logical reasoning optimization, pushing the 525 boundaries of what LLMs can achieve in structured 526 reasoning tasks. However, RL enhances sampling 527 efficiency of existing reasoning paths encoded 528 in the base model but does not generate new reasoning patterns (Yue et al., 2025) and therefore 530 relies on the capabilities of the foundation models. From this perspective, SFT with high-quality 532 human-annotated or LLM-distilled data is a simpler and more effective way to expand the 534 boundaries of smaller models than RL.

5.2.3 Inference-Time Decoding

537

541

542

545

546

547

548

551

552

553

554

556

558

560

561

565

566

We categorize logical reasoning enhancement methods during inference-time into inference-time scaling and constrained decoding.

Inference-time scaling employs computational augmentation without parameter updates. One common approach is decoding with structured outputs and modular workflows. GoT (Lei et al., 2023) creates structured reasoning nodes to improve complex multi-step logical reasoning. Similarly, Chain of Logic (Servantez et al., 2024) introduces a method that first divides the legal question into smaller sub-questions for solving, and then assembles the answers to resolve the original problem. In other contexts, researchers design more complex modular workflows for better performance (Creswell et al., 2023; Malon et al., 2024).

Another inference-time scaling approach involves stimulating autonomous reasoning, guiding LLMs to iteratively refine their answers. Maieutic Prompting (Jung et al., 2022) eliminates contradictions through recursive reasoning. Similarly, Logicof-Thoughts (Liu et al., 2024a) and DetermLR (Sun et al., 2024) progressively approach the answers in an iterative style.

Inference-time scaling offers the flexibility to improve model performance without additional training, but usually incurs higher inference costs.

Constrained decoding methods, on the other hand, focus on improving the controllability and reliability of reasoning processes. Neurologic (Lu et al., 2021) enforces predicate logic constraints, while Formal-LLM (Li et al., 2024) integrates automata for constraining plan generation.

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

5.3 External Knowledge Utilization

LLMs often generate incorrect answers due to hallucinations when performing complex tasks such as logical reasoning, making it necessary to incorporate external knowledge to assist in producing accurate responses. Formally, the optimal integration of external knowledge can be formulated as a joint optimization problem:

$$(M^*, K^*) = \arg \max_{M, K} R(M, K)$$
(3)

where:

- *M*: the neural model, which includes both the model's parameters and its decoding strategies (generally with the parameters unchanged).
- *K*: knowledge integration strategy, including knowledge source curation, structured knowledge representation, retrieval-augmented mechanisms, etc.
- *R*: reasoning performance evaluator.

Zayyad and Adi (2024) and Yang et al. (2023) extract data from Lean, a mathematical proof tool, to aid theorem proving. In contrast, "Logic-Query-of-Thoughts" (LQOT) (Liu et al., 2024b) decomposes complex logical problems into easier sub-questions before integrating knowledge graphs.

In reading comprehension, Ouyang et al. (2023) construct supergraphs to address complex contextual reasoning, while KnowRA (Mai et al., 2025) autonomously determines whether to accept external knowledge to assist document-level relation extraction.

5.4 Neuro-Symbolic Approaches

Neural-symbolic hybrid methods represent a burgeoning research area that aims to combine the powerful representational capabilities of deep learning with the precision and interpretability of symbolic reasoning.

Formally, a neural-symbolic hybrid system aims to optimize both the neural model M and the symbolic solver P (where P represents the symbolic reasoning process) to maximize logical reasoning performance. The overall objective can be expressed as:

614

615

617

618

622

623

624

625

631

632

634

638

639

643

647

651

652

655

where:

 M: The neural model, which includes both the model's parameters and its decoding strategies. It maps the input x (e.g., natural language) into a symbolic representation z within a formal language L:

 $(M^*, P^*) = \arg\max_{M, P} R(P(M(x))),$

$$z = M(x), \quad z \in \mathcal{L}.$$

• *P*: The symbolic solver, which operates on the symbolic representation *z* produced by *M* to generate the final output *y*:

$$y = P(z).$$

• *R*: The reasoning performance metric.

Two key directions of the optimization process:

- Improving *M*: including refining the model's parameters and decoding strategies to produce symbolic representations that are both accurate and compatible with *P*.
- Enhancing *P*: involving improving the symbolic solver's ability to process.

By jointly optimizing M and P, neural-symbolic hybrid systems aim to leverage the strengths of both neural networks and symbolic reasoning to achieve superior logical reasoning capabilities. It is worth noting that in earlier neural-symbolic pipelines, Pis often implemented as a fixed external logical reasoning engine, and thus is generally not optimized. However, in advanced practice, LLMs are increasingly being used to perform the role of P, enabling diverse optimization.

Fundamentally, these methods involve translating problems into symbolic representations with LLMs, and external symbolic solvers solving them. For example, in LINC (Olausson et al., 2023), LLMs convert natural language (NL) into firstorder logic (FOL) expressions, and utilize an external theorem prover for deductive inference.

Further efforts focus on improving NL-tosymbolic translation. One prevailing approach is directly optimizing translation through training (Yang et al., 2024a) or decoding strategies (Ryu et al., 2024), while the others depend on verification or correction mechanisms (Yang et al., 2024a; Pan et al., 2023).

Building upon these, recent advancements address the traditional pipeline limitations by fully integrating LLMs into reasoning processes. Logic Agent (LA) (Liu et al., 2024a) replaces external solvers with rule-guided LLM inference chains, while LLM-TRes (Toroghi et al., 2024) implements self-contained verifiable reasoning without external symbolic solvers. SymbCoT (Xu et al., 2024c) coordinates translation, planning, solving and verification entirely through LLMs. Xu et al. (2024b) propose Aristotle, which further systematizes the symbolic reasoning pipeline through three LLMdriven components: Logical Decomposer, Logical Search Router, and Logical Resolver. However, these modular approaches increase system complexity and do not fundamentally enhance the models' intrinsic reasoning capabilities.

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

6 Discussion

The landscape of logical reasoning in LLMs presents several unresolved challenges that merit deeper examination. In Appendix A.2, we provide an extensive discussion on fundamental tensions shaping current research: the gap between surface-level pattern matching and genuine logical competence; inconsistencies in robustness across varied reasoning contexts; competing priorities of interpretability versus performance in hybrid approaches; limitations in current evaluation paradigms; future directions.

7 Conclusion

This survey synthesizes the rapid advancements and persistent challenges in logical reasoning for large language models (LLMs). While LLMs demonstrate impressive heuristic reasoning, rigorous logical inference remains inconsistent due to limitations in robustness, generalization, and interpretability. We analyzed strategies to enhance reasoning, including neuro-symbolic integration, data-centric tuning, reinforcement learning, testtime scaling and other improved decoding methods, and highlighted benchmarks like FOLIO and LogiQA for systematic evaluation. Future progress hinges on hybrid architectures that unify neural and symbolic reasoning, robust evaluation frameworks, scalable methods for cross-domain and multimodal inference, and directly enhancing base model causality. Addressing these challenges will advance LLMs toward reliable, interpretable reasoning critical for real-world applications.

712

714

715

716

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

734

736

737

739

740

741

742

744

745

746

747

748

749

750

751

753

5 Limitations

We focus on formal logical reasoning, which is in
line with the symbolic approaches actively studied
in the literature. For general reasoning or other
specific reasoning types, readers may refer to complementary surveys.

711 References

- Guangsheng Bao, Hongbo Zhang, Cunxiang Wang, Linyi Yang, and Yue Zhang. 2025. How likely do LLMs with CoT mimic human reasoning? In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7831–7850, Abu Dhabi, UAE. Association for Computational Linguistics.
 - Andrew M Bean, Simi Hellsten, Harry Mayne, Jabez Magomere, Ethan A Chi, and 1 others. 2024. Lingoly: A benchmark of olympiad-level linguistic reasoning puzzles in low-resource and extinct languages. *arXiv preprint arXiv:2406.06196*.
 - Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen tau Yih, and Yejin Choi. 2020. Abductive commonsense reasoning. In International Conference on Learning Representations.
- Chen Bowen, Rune Sætre, and Yusuke Miyao. 2024. A comprehensive evaluation of inductive reasoning capabilities and problem solving in large language models. In *Proc. of ACL Findings*, pages 323–339.
- Ronnie Cann. 1993. *Formal semantics: an introduction*. Cambridge University Press, United States.
- Nuri Cingillioglu and Alessandra Russo. 2019. Deeplogic: Towards end-to-end differentiable logical reasoning. *Preprint*, arXiv:1805.07433.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2021. Transformers as soft reasoners over language. In *Proc. of IJCAI*.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2023. Selection-inference: Exploiting large language models for interpretable logical reasoning. In *Proc.* of *ICLR*.
- DeepSeek-AI. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. Technical report.
- Maksym Del and Mark Fishel. 2023. True detective: A deep abductive reasoning benchmark undoable for GPT-3 and challenging for GPT-4. In *Proceedings* of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023), pages 314–322.

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, and 1 others. 2023. Faith and fate: Limits of transformers on compositionality. *arXiv preprint arXiv:2305.18654*.

754

755

758

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

- Jiazhan Feng, Ruochen Xu, Junheng Hao, Hiteshi Sharma, Yelong Shen, and 1 others. 2024. Language models can be deductive solvers. In *Proc. of ACL Findings*, pages 4026–4042.
- João Pedro Gandarela, Danilo S Carvalho, and André Freitas. 2024. Inductive learning of logical theories with llms: A complexity-graded analysis. *arXiv preprint arXiv:2408.16779*.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2022. ROSCOE: A suite of metrics for scoring step-by-step reasoning.
- Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, and 1 others. 2024a. FOLIO: Natural language reasoning with first-order logic. In *Proc. of EMNLP*, pages 22017–22031.
- Simeng Han, Aaron Yu, Rui Shen, Zhenting Qi, Martin Riddell, and 1 others. 2024b. P-FOLIO: Evaluating and improving logical reasoning with abundant human-written reasoning chains. In *Proc. of EMNLP Findings*, pages 16553–16565.
- Izumi Haruta, Koji Mineshima, and Daisuke Bekki. 2020. Logical inferences with comparatives and generalized quantifiers. In *Proc. of ACL*, pages 263–270.
- Chadi Helwe, Chloé Clavel, and Fabian Suchanek. 2022. Logitorch: A pytorch-based library for logical reasoning on natural language. In *Proc. of EMNLP*.
- Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.
- Jin Jiang, Yuchen Yan, Yang Liu, Yonggang Jin, Shuai Peng, and 1 others. 2024. Logicpro: Improving complex logical reasoning via program-guided learning. *arXiv preprint arXiv:2409.12929*.
- Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F. Chen, and Shafiq Joty. 2024. Learning planningbased reasoning by trajectories collection and process reward synthesizing. In *Proc. of EMNLP*, pages 334– 350.
- Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, and 1 others. 2022. Maieutic prompting: Logically consistent reasoning with recursive explanations. In *Proc. of EMNLP*, pages 1266–1279.

910

911

912

913

914

Bin Lei, Chunhua Liao, Caiwen Ding, and 1 others. 2023. Boosting logical reasoning in large language models through a new framework: The graph of thought. *arXiv preprint arXiv:2308.08614*.

810

811

812

813

814

815

817

819

822

825

829

832

834

835

836

837

839

845

847

848

850

851

853 854

855

856

- Jacqueline P. Leighton. 2003. *Defining and Describing Reason*. Cambridge University Press, Cambridge, UK.
- Zelong Li, Wenyue Hua, Hao Wang, He Zhu, and Yongfeng Zhang. 2024. Formal-Ilm: Integrating formal language and natural language for controllable Ilm-based agents. *arXiv preprint arXiv:2402.00798*.
- Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong, Zhijiang Guo, Le Song, and Cheng-Lin Liu. 2025. From system 1 to system 2: A survey of reasoning large language models. *Preprint*, arXiv:2502.17419.
 - Hanmeng Liu, Leyang Cui, Jian Liu, and Yue Zhang. 2021a. Natural language inference in context - investigating contextual reasoning over long texts. *Proc.* of AAAI, pages 13388–13396.
 - Hanmeng Liu, Yiran Ding, Zhizhang Fu, Chaoli Zhang, Xiaozhang Liu, and Yue Zhang. 2025. Evaluating the logical reasoning abilities of large reasoning models. *Preprint*, arXiv:2505.11854.
 - Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, and 1 others. 2023a. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, pages 2947–2962.
 - Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023b. Evaluating the logical reasoning ability of chatgpt and gpt-4. *Preprint*, arXiv:2304.03439.
 - Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. 2023c. Logicot: Logical chain-of-thought instruction tuning. In *Proc. of EMNLP Findings*, pages 2908–2921.
 - Hanmeng liu, Zhiyang Teng, Ruoxi Ning, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Glore: Evaluating logical reasoning of large language models. *Preprint*, arXiv:2310.09107.
 - Hanmeng Liu, Zhiyang Teng, Chaoli Zhang, and Yue Zhang. 2024a. Logic agent: Enhancing validity with logic rule invocation. *Preprint*, arXiv:2404.18130.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021b. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning.
- Lihui Liu, Zihao Wang, Ruizhong Qiu, Yikun Ban, Eunice Chan, and 1 others. 2024b. Logic query of thoughts: Guiding large language models to answer complex logic queries with knowledge graphs. *Preprint*, arXiv:2404.04264.

- Yinhong Liu, Zhijiang Guo, Tianya Liang, Ehsan Shareghi, Ivan Vulić, and Nigel Collier. 2024c. Aligning with logic: Measuring, evaluating and improving logical consistency in large language models. *arXiv preprint arXiv:2410.02205.*
- Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. NeuroLogic decoding: (un)supervised neural text generation with predicate logic constraints. In *Proc. of NAACL*, pages 4288–4299.
- Man Luo, Shrinidhi Kumbhar, Ming shen, Mihir Parmar, Neeraj Varshney, and 1 others. 2024. Towards logiglue: A brief survey and a benchmark for analyzing logical reasoning capabilities of language models. *Preprint*, arXiv:2310.00836.
- Chengcheng Mai, Yuxiang Wang, Ziyu Gong, Hanxiang Wang, and Yihua Huang. 2025. Knowra: Knowledge retrieval augmented method for document-level relation extraction with comprehensive reasoning abilities. *Preprint*, arXiv:2501.00571.
- Christopher Malon, Martin Min, Xiaodan Zhu, and 1 others. 2024. Exploring the role of reasoning structures for constructing proofs in multi-step natural language reasoning with large language models. In *Proc. of EMNLP*, pages 15299–15312.
- J. McCarthy and P.J. Hayes. 1981. Some philosophical problems from the standpoint of artificial intelligence. In *Readings in Artificial Intelligence*, pages 431–450.
- John McCarthy. 1959. Programs with common sense. In Proceedings of the Teddington Conference on the Mechanization of Thought Processes.
- John McCarthy. 1989. Artificial intelligence, logic and formalizing common sense. *Philosophical Logic and Artificial Intelligence*, pages 161–190.
- Philipp Mondorf and Barbara Plank. 2024. Beyond accuracy: Evaluating the reasoning behavior of large language models–a survey. *arXiv preprint arXiv:2404.01869*.
- Terufumi Morishita, Gaku Morio, Atsuki Yamaguchi, and Yasuhiro Sogawa. 2024. Enhancing reasoning capabilities of Ilms via principled synthetic logic corpus. In *Proc. of NeurIPS*, pages 73572–73604.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, and 1 others. 2025. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*.
- A. Newell and H. Simon. 1956. The logic theory machine–a complex information processing system. *IRE Transactions on Information Theory*.
- Ha-Thanh Nguyen, Randy Goebel, Francesca Toni, Kostas Stathis, and Ken Satoh. 2023. How well do sota legal reasoning models support abductive reasoning? *Preprint*, arXiv:2304.06912.

915

963 964

- 967
- 965 966

969

- Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, and 1 others. 2023. LINC: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In Proc. of EMNLP, pages 5153-5176.
- OpenAI. 2024. Learning to reason with LLMs. Technical report.
- Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. 2023. Fact-driven logical reasoning for machine reading comprehension. Preprint, arXiv:2105.10334.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning. In Proc. of EMNLP Findings, pages 3806-3824.
- Mihir Parmar, Nisarg Patel, Neeraj Varshney, Mutsumi Nakamura, Man Luo, and 1 others. 2024. Logicbench: Towards systematic evaluation of logical reasoning ability of large language models. In Proc. of ACL, pages 13679–13707.
- Mihir Parmar, Neeraj Varshney, Nisarg Patel, Santosh Mashetty, Man Luo, and 1 others. 2023. Logicbench: A benchmark for evaluation of logical reasoning.
- Fernando Carlos Neves Pereira. 1982. Logic for natural language analysis.
- Molly Petersen and Lonneke van der Plas. 2023. Can language models learn analogical reasoning? investigating training objectives and comparisons to human performance. In Proc. of EMNLP, pages 16414-16425.
 - Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. 2023. ReCEval: Evaluating reasoning chains via correctness and informativeness. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 10066-10086, Singapore. Association for Computational Linguistics.
- Chengwei Qin, Wenhan Xia, Tan Wang, Fangkai Jiao, Yuchen Hu, and 1 others. 2024. Relevant or random: Can llms truly perform analogical reasoning? Preprint, arXiv:2404.12728.
- Kanagasabai Rajaraman, Saravanan Rajamanickam, and Wei Shi. 2023. Investigating transformer-guided chaining for interpretable natural logic reasoning. In Proc. of ACL Findings, pages 9240-9253.
- Samuel Ryb, Mario Giulianelli, Arabella Sinclair, and Raquel Fernández. 2022. AnaLog: Testing analytical and deductive logic learnability in language models. In Proceedings of the 11th Joint Conference on Lexical and Computational Semantics, pages 55-68.
- Hyun Ryu, Gyeongman Kim, Hyemin S Lee, and Eunho Yang. 2024. Divide and translate: Compositional first-order logic translation and verification for complex logical reasoning. arXiv preprint arXiv:2410.08047.

Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi, and 1 others. 2023. Testing the general deductive reasoning capacity of large language models using ood examples. In Proc. of NeurIPS, pages 3083–3105.

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1002

1003

1004

1005

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

- Sergio Servantez, Joe Barrow, Kristian Hammond, and Rajiv Jain. 2024. Chain of logic: Rule-based reasoning with large language models. In Proc. of ACL Findings, pages 2721–2733.
- Jihao Shi, Xiao Ding, Li Du, Ting Liu, and Bing Qin. 2021. Neural natural logic inference for interpretable question answering. In Proc. of EMNLP, pages 3673-3684.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. arXiv preprint arXiv:2305.17493.
- Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. 2019. Clutrr: A diagnostic benchmark for inductive reasoning from text. Empirical Methods of Natural Language Processing (EMNLP).
- Yueqi Song, Tianyue Ou, Yibo Kong, Zecheng Li, Graham Neubig, and Xiang Yue. 2025. Visualpuzzles: Decoupling multimodal reasoning evaluation from domain knowledge. Preprint, arXiv:2504.10342.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, and 1 others. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. arXiv preprint arXiv:2206.04615.
- Joe Stacey, Pasquale Minervini, Haim Dubossarsky, and Marek Rei. 2022. Logical reasoning with span-level predictions for interpretable and robust NLI models. In Proc. of EMNLP, pages 3809–3823.
- Michael Sullivan. 2024. It is not true that transformers are inductive learners: Probing NLI models with external negation. In Proc. of EACL, pages 1924-1945.
- Hongda Sun, Weikai Xu, Wei Liu, Jian Luan, Bin Wang, and 1 others. 2024. Determlr: Augmenting llm-based logical reasoning from indeterminacy to determinacy. In Proc. of ACL, pages 9828–9862.
- Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, and 1 others. 2023. A survey of reasoning with foundation models. arXiv preprint arXiv:2312.11562.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. 1018 ProofWriter: Generating implications, proofs, and 1019 abductive statements over natural language. In Proc. 1020 of ACL Findings, pages 3621-3634. 1021

Ramya Keerthy Thatikonda, Wray Buntine, and Ehsan Jundong Xu, Hao Fei, Meng Luo, Qian Liu, Liang-1075 Shareghi. 2025. Assessing the alignment of fol closeming Pan, and 1 others. 2024b. Aristotle: Mas-1076 ness metrics with human judgement. arXiv preprint tering logical reasoning with a logic-complete arXiv:2501.08613. decompose-search-resolve framework. arXiv preprint arXiv:2412.16953. 1079 Jidong Tian, Yitian Li, Wenqing Chen, Liqiang Xiao, Hao He, and Yaohui Jin. 2021. Diagnosing the first-Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-1080 order logical reasoning ability through LogicNLI. In Li Lee, and Wynne Hsu. 2024c. Faithful logical 1081 Proc. of EMNLP, pages 3738-3747. reasoning via symbolic chain-of-thought. In Proc. of 1082 ACL, pages 13326–13365. 1083 Armin Toroghi, Willis Guo, Ali Pesaranghader, and Scott Sanner. 2024. Verifiable, debuggable, and Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chala-1084 repairable commonsense logical reasoning via llmmala, Peiyang Song, and 1 others. 2023. Leandojo: 1085 based theory resolution. In Proc. of EMNLP, pages theorem proving with retrieval-augmented language 1086 6634-6652. models. In Proc. of ICONIP. 1087 Chen Wang, Xudong Li, Haoran Liu, Xinyue Wu, and Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, 1088 Wanting He. 2024a. Efficient logical reasoning in and Faramarz Fekri. 2024a. Harnessing the power of 1089 large language models through program-guided learnlarge language models for natural language to first-1090 ing. Authorea Preprints. order logic translation. In Proc. of ACL, pages 6942-6959. 1092 Chenxu Wang, Ping Jian, and Zhen Yang. 2024b. Thought-path contrastive learning via premise-Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik oriented data augmentation for logical reading com-1093 prehension. arXiv preprint arXiv:2409.14495. Cambria, and 1 others. 2024b. Language models as 1094 inductive reasoners. In Proc. of EACL, pages 209-1095 Po-Wei Wang, Priya L. Donti, Bryan Wilder, and Zico 225. 1096 Kolter. 2019. Satnet: Bridging deep learning and logical reasoning using a differentiable satisfiability Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 1097 solver. Preprint, arXiv:1905.12149. 2020. Reclor: A reading comprehension dataset re-1098 quiring logical reasoning. In Proc. of ICLR. 1099 Siyuan Wang, Zhongkun Liu, Wanjun Zhong, Ming Zhou, Zhongyu Wei, and 1 others. 2022. From lsat: Zhangdie Yuan, Songbo Hu, Ivan Vulić, Anna Korho-1100 The progress and challenges of complex reasoning. nen, and Zaiqiao Meng. 2023. Can pretrained lan-1101 IEEE/ACM Transactions on Audio, Speech, and Language models (yet) reason deductively? In Proc. of 1102 guage Processing. EACL, pages 1447-1462. 1103 Thilini Wijesiriwardene, Ruwan Wickramarachchi, Bi-Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai 1104 mal Gajera, Shreeyash Gowaikar, Chandan Gupta, Wang, Shiji Song, and Gao Huang. 2025. Does reand 1 others. 2023. ANALOGICAL - a novel bench-1105 inforcement learning really incentivize reasoning ca-1106 mark for long text analogy evaluation in large language models. In Proc. of ACL Findings, pages 3534pacity in llms beyond the base model? arXiv preprint 1107 3549. arXiv:2504.13837. 1108 Zhiheng Xi, Wenxiang Chen, Boyang Hong, Senjie Majd Zayyad and Yossi Adi. 2024. Formal language 1109 Jin, Rui Zheng, and 1 others. 2024. Training large knowledge corpus for retrieval augmented generation. 1110 language models for reasoning through reverse cur-Preprint, arXiv:2412.16689. 1111 riculum reinforcement learning. In Proc. of ICML. Matej Zečević, Moritz Willig, Devendra Singh Dhami, 1112 Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. and Kristian Kersting. 2023. Causal parrots: Large 1113 2024. Logicvista: Multimodal llm logical realanguage models may talk causality but are not 1114 soning benchmark in visual contexts. Preprint, causal. Transactions on Machine Learning Research, 1115 arXiv:2407.04973. 2023(8). 1116 Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiang-1117 Jun Liu, and Erik Cambria. 2023. Are large lanming Shu, Jinlin Xiao, and 1 others. 2024. o1-1118 guage models really good logical reasoners? a coder: an o1 replication for coding. arXiv preprint 1119 comprehensive evaluation and beyond. Preprint, arXiv:2412.00154. 1120 arXiv:2306.09841. Fangzhi Xu, Zhiyong Wu, Qiushi Sun, Siyu Ren, Fei Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi 1121 Yuan, and 1 others. 2024a. Symbol-LLM: Towards Shi, and 1 others. 2024. Marco-o1: Towards open 1122 foundational symbol-centric interface for large lanreasoning models for open-ended solutions. arXiv 1123 guage models. In Proc. of ACL, pages 13091-13116. preprint arXiv:2411.14405. 1124

1022

1023

1024

1026

1030

1032

1033

1035

1037

1039 1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056 1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

A Appendix

A.1 Data Examples

Figure 2 presents example questions from datasets related to logical reasoning, serving as a supplement to the description of task types in the main text.

A.2 Further Discussion

The integration of logical reasoning into large language models (LLMs) remains a critical challenge, marked by persistent gaps between heuristic performance and formal logical rigor. Below, we analyze three unresolved tensions dominating the field and outline future directions.

Superficial Reasoning vs. Genuine Logical Com-1138 petence. Despite promising results on benchmark 1139 datasets, an fundamental question persists: do 1140 LLMs truly possess logical reasoning abilities, or 1141 do they merely approximate reasoning through so-1142 phisticated pattern recognition? Some recent stud-1143 ies have observed that current LLMs often lack sub-1144 stantive causal reasoning (Bao et al., 2025; Zečević 1145 et al., 2023), which in turn limits both their capacity 1146 for genuine inference and their ability to generalize, 1147 posing a crucial challenge to be tackled. 1148

Robustness vs. Generalization. LLMs exhibit 1149 1150 inconsistent performance in structured reasoning tasks such as deductive inference and abductive 1151 hypothesis generation. While models fine-tuned on 1152 datasets like FOLIO (Han et al., 2024a) excel in 1153 controlled settings, they struggle with adversarial 1154 perturbations or semantically equivalent rephras-1155 ings. This inconsistency arises from their reliance 1156 on surface-level statistical correlations rather than 1157 causal relationships, coupled with limited out-of-1158 distribution generalization. A key question per-1159 sists: can LLMs achieve human-like robustness 1160 without sacrificing cross-domain adaptability? Cur-1161 rent methods prioritize narrow task performance, 1162 leaving real-world applicability uncertain. 1163

Interpretability vs. Performance. A central ten-1164 sion lies in balancing neural scalability with sym-1165 bolic precision. Neuro-symbolic approaches like 1166 Logic-LM (Pan et al., 2023) and Symbol-LLM (Xu 1167 1168 et al., 2024a) embed formal logic solvers into neural architectures, improving interpretability through 1169 step-by-step proofs. However, these methods face 1170 scalability bottlenecks with large knowledge bases 1171 or complex rule dependencies. Conversely, data-1172



(a) A multi-choice reading comprehension example from the LogiQA dataset.



(b) An NLI example from the ConTRoL dataset.

Figure 2: Example tests of Logical reasoning in NLP tasks.

driven methods (e.g., instruction tuning on LogicBench (Parmar et al., 2024)) achieve broader task coverage but fail to generalize beyond syntactic patterns. How can we reconcile transparent reasoning with black-box model performance? Hybrid architectures offer promise but introduce computational overhead, limiting practical deployment.

1173

1174

1175

1176

1177

1178

1179

Evaluation Rigor. Existing benchmarks like 1180 LogiQA (Liu et al., 2021b) and ReClor (Yu et al., 1181 2020) conflate reasoning ability with pattern recog-1182 nition through multiple-choice formats. While ef-1183 forts like NeuLR (Xu et al., 2023) curate "neutral" 1184 content to isolate reasoning from domain knowl-1185 edge, they lack scope for holistic evaluation. Cur-1186 rent metrics (e.g., accuracy, BLEU) fail to assess 1187 consistency (invariance to logically equivalent in-1188 puts) or soundness (adherence to formal proof struc-1189 tures). What defines a gold standard for logical 1190 reasoning evaluation? Benchmarks must prioritize 1191 systematic testing of core principles (e.g., transitiv-1192 ity, contraposition) over task-specific performance. 1193

1194 Future Directions. Addressing these challenges requires hybrid architectures that dynamically in-1195 tegrate neural and symbolic components, such as 1196 differentiable theorem provers, to balance scala-1197 bility and precision. Equally important is the de-1198 velopment of evaluation frameworks that stress-1199 test models on perturbed logical statements (e.g., 1200 negated premises, swapped quantifiers) to isolate 1201 reasoning from memorization. Multi-modal reason-1202 ing, which grounds inference in diverse modalities 1203 (text, images, code), presents untapped potential 1204 for enhancing robustness and interpretability. Fi-1205 nally, interdisciplinary collaboration-leveraging 1206 insights from formal logic, cognitive science, and 1207 machine learning-will be essential to design sys-1208 tems that reason with and about uncertainty. Until 1209 LLMs reliably disentangle logic from lexicon, their 1210 deployment in high-stakes domains will remain 1211 precarious. Bridging this gap demands rigorous 1212 benchmarks, scalable hybrid methods, and a redefi-1213 nition of evaluation paradigms. 1214