

Logical Reasoning over Natural Language as Knowledge Representation: A Survey

Anonymous ACL submission

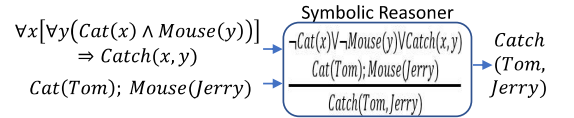
Abstract

Logical reasoning is central to human cognition and intelligence. Past research of logical reasoning within AI uses formal language as knowledge representation (and symbolic reasoners). However, reasoning with formal language has proved challenging (e.g., brittleness and knowledge-acquisition bottleneck). This paper provides a comprehensive overview on a new paradigm of logical reasoning, which uses natural language as knowledge representation (and pretrained language models as reasoners), including philosophical definition and categorization of logical reasoning, advantages of the new paradigm, benchmarks and methods, challenges of the new paradigm, possible future directions, and relation to related NLP fields. This new paradigm is promising since it not only alleviates many challenges of formal representation but also has advantages over end-to-end neural methods.

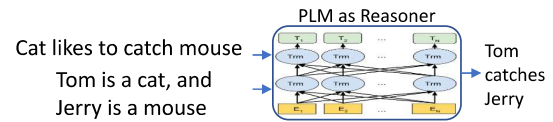
1 Introduction

An argument consists of premise(s) and a conclusion. Logical reasoning is a form of thinking in which premises and relations between premises are used in a rigorous manner to infer conclusions that are entailed (or implied) by the premises and the relations (Nunes, 2012). It consists of three reasoning types, namely deductive reasoning, inductive reasoning, and abductive reasoning (Flach and Kakas, 2000) (more illustration on the categorization can be found in §2). It is important since the ability to reach logical conclusions on the basis of prior information is recognized as central to human cognition and intelligence (Goel et al., 2017).

The past research of logical reasoning within AI uses formal language (e.g., first-order logic) as knowledge representation and symbolic reasoners (Muggleton and Raedt, 1994). This paradigm has resulted in impressive applications such as expert systems (Metaxiotis et al., 2002). However, building and reasoning over formal language



(a) Formal language as knowledge representation and symbolic reasoner



(b) Natural language as knowledge representation and PLM as reasoner

Figure 1: Comparison between the previous paradigm which uses formal representation and symbolic reasoner, and the new paradigm which uses natural language as knowledge representation and PLM as reasoner.

have proved challenging (Musen and Van der Lei, 1988), with representative disadvantages of brittleness (an expert system fails as long as its knowledge base does not contain complete knowledge for a problem) and knowledge-acquisition bottleneck (human experts are needed to encode their knowledge with formal representation).

Since the rapid development in language models, natural language has been explored as a new knowledge representation, and pretrained language model (PLM) has been used as a new corresponding reasoner for deductive reasoning (Clark et al., 2020), abductive reasoning (Bhagavatula et al., 2020), and inductive reasoning (Yang et al., 2022b). Therefore, all three reasoning types of logical reasoning have been investigated with natural language as knowledge representation. These research also shows that PLMs can be finetuned or prompted to perform well for each of the reasoning types.

In this paper, we summarize the three previously separately investigated logical reasoning types together, referred as logical reasoning from the perspectives of deductive, inductive, and abductive reasoning over natural language as knowledge representation (LRNL), and provide an in-depth and up-to-date survey of LRNL.

Illustrated in Figure 1, LRNL means a new paradigm for logical reasoning that uses new knowledge representation (natural language) and new reasoner (PLM). Recent methods in this area are generally modular-based: multiple PLMs each as one module playing a different function, combined together to perform complex tasks. They make one step of reasoning with one inference of PLM. For complex problems, they usually have access to a knowledge base that stores relevant textual knowledge to be retrieved as premises to support the reasoning process to reach a conclusion, which might be used as a new premise for the next step’s reasoning. By iteratively repeating this process, a final conclusion may be made. Although looks similar to expert systems, we discuss how LRNL is possible to overcome many main challenges of the previous paradigm such as brittleness and knowledge-acquisition bottleneck in §3.1.

In addition to the comparison with formal language, in §3.2 we discuss that LRNL could be viewed as a new type of neural-symbolic (NeSy) method, which has unique advantages over existing NeSy methods. We also discuss how LRNL, as a NeSy method, has advantages over existing end-to-end neural methods (e.g., explainability, controllability, less catastrophic forgetting) in §3.3. These advantages make an LRNL system possible to deal with many challenging problems today.

In the remaining sections of this survey, we review papers on LRNL (including deductive reasoning §4, inductive reasoning §5, and abductive reasoning §6), and list challenges (§7 and §A.9). Our main focus is to understand the language model’s logical reasoning ability through the three sub-types of logical reasoning to provide finer analysis and avoid ambiguity. Therefore we focus on papers that specialize on one (or more) of the three sub-types of logical reasoning (instead of only “reasoning”). In §A.1 we discuss the relation of LRNL to related NLP fields (e.g., commonsense reasoning), which could help to form a clear shape of LRNL in NLP. For each reasoning sub-type, we summarize existing task formulations, datasets, and methods under each task.

2 Definition and Categorization of Logical Reasoning

There are many subjects related to logical reasoning, including philosophy, logic, and AI. Among them, the definition and categorization aspects of

logical reasoning are handled by philosophy research. However, debate exists in philosophy research on the categorization of logical reasoning. We leave a detailed description of the debate in philosophy research in § A.2 and only leave the conclusions here according to philosophy research.

In general, logical reasoning consists of deductive, inductive, and abductive reasoning (Console and Saitta, 2000). Given an argument consisting of premises and a conclusion, we define the sub-type of logical reasoning it involves below:

Definition for deductive reasoning: the premises can conclusively provide support for the conclusion, i.e. if the premises are all true, it would be impossible for the conclusion to be false.

Definition for inductive reasoning: the premises cannot conclusively provide support for the conclusion, since the conclusion generalizes existing information in premises to new knowledge, which has a wider applicable scope than those in premises.

Definition for abductive reasoning: the premises cannot conclusively provide support for the conclusion, since the conclusion contains more specific information over the premises (most commonly used as generating most probable explanations).

Please note that according to Console and Saitta (2000), inductive reasoning and abductive reasoning are not exclusive to each other.

3 Advantages of LRNL

3.1 Advantages over Formal Language

Building and reasoning over formal language have proved challenging (Musen and Van der Lei, 1988; Cropper et al., 2022), with disadvantages such as (1) brittleness (expert system fails when its knowledge base does not contain complete knowledge for a problem), (2) knowledge-acquisition bottleneck (human experts are needed to encode their knowledge with formal representation), (3) inability to handle raw data such as natural language, (4) sensitivity to label errors, and (5) failure to recognize different symbols with similar meanings.

Nevertheless, the new paradigm of logical reasoning, LRNL, has systematic strengths over these challenges. Specifically, PLMs contain knowledge themselves (Davison et al., 2019), which makes it possible for them to provide good answers even when some required explicit knowledge is not present in a knowledge base (Talmor et al., 2020) (less brittle), and be less affected by input errors (Meng et al., 2021); with natural language as

Dataset	Human written	Realistic	Multi-step	Theory included	Theory sufficient	Proof generation	Size
D*	✗	✗	✓	✓	✓	✗	500k
ParaRules	✓	✗	✓	✓	✓	✗	40k
Birds-electricity	✓	✓	✓	✓	✓	✗	5k
Leap-of-thought	✗	✓	✗	✓	✗	✗	33k
PARARULE-Plus	✗	✗	✓	✓	✓	✗	400k
FOLIO	✓	✓	✓	✓	✓	✗	1,435
D*(CWA)	✗	✗	✓	✓	✓	✓	500k
D*(OWA)	✗	✗	✓	✓	✗	✓	500k
EntailmentBank	✓	✓	✓	✓	✓	✓	1,840
ENWN	✓	✓	✓	✓	✓	✓	100

Table 1: Summary of deductive reasoning datasets: D*, ParaRules, and birds-electricity (Clark et al., 2020); leap-of-thought (Talmor et al., 2020); PARARULE-Plus (Bao et al., 2022); FOLIO (Han et al., 2022); D*(CWA) and D*(OWA) (Tafjord et al., 2021); EntailmentBank (Dalvi et al., 2021); ENWN (Sprague et al., 2022).

knowledge representation, such a system can naturally handle raw input, and is possible to utilize the enormous web corpora to automatically construct a rule base using information extraction (Ji, 2018) or inductive reasoning (Yang et al., 2022b) (less affected by knowledge-acquisition bottleneck); using embeddings for concepts (Mikolov et al., 2013), it semantically “understands” the meaning of symbols and therefore robust for paraphrasing.

3.2 Advantages over Existing NeSy Systems

LRNL could be seen as a new type of NeSy in addition to the existing 6 types summarized by Kautz (2022), as its goal and design of methodology are typically symbolic (logical reasoning with knowledge bases), while avoiding any symbolic representation, using (currently pure) neural methods. Therefore LRNL can avoid many bottlenecks of the other NeSy methods caused by symbolic representation, such as symbolic knowledge acquisition and scalability (Wang and Yang, 2022).

3.3 Advantages over E2E Neural Methods

As a NeSy method, LRNL systematically has some advantages over end-to-end neural methods, such as interpretability (Cambria et al., 2023) (since its stepwise reasoning nature), more controllability (LRNL reasons following a given knowledge base), and less catastrophic forgetting (LRNL uses an explicit knowledge base to store knowledge).

4 Deductive Reasoning

4.1 Existing Task Formulations

Existing tasks for deductive reasoning can be summarized as hypothesis classification, proof generation, proof generation with incomplete information,

and implication enumeration. Datasets for tasks are summarized in Table 1. “Proof generation” tab with ✗ means it is for hypothesis classification task.

Hypothesis Classification Each data example for hypothesis classification task is a tuple $(theory, hypothesis, correctness)$, where $theory$ typically has the form $(fact^*, rule^*)$, $hypothesis$ is a question, and $correctness$ can be *True* or *False* (or *Unknown*). This task requires to predict the $correctness$ for the $hypothesis$ given the $theory$.

Proof Generation The proof generation task has the same setting as the hypothesis classification task, except that in addition to predicting a $correctness$, the proof generation task also requires providing a $proof$ given $theory$ to explain the $correctness$. The $proof$ is a directed tree $(\mathcal{N}, \mathcal{E})$ with nodes $n \in \mathcal{N}$ and edges $e \in \mathcal{E}$. Each node is an item of knowledge in $theory$ (usually a $fact$ or a $rule$), or a generated intermediate reasoning conclusion, or the $hypothesis$ itself; Each edge points from a premise node to a conclusion node to form a deductive argument, which typically needs one-step inference (not multi-step).

Proof Generation with Incomplete Information This task is the same as the proof generation task, except that $theory$ lacks one $node$ to form a complete $proof$. Specifically, given $theory$, it requires to predict the $correctness$ of $hypothesis$ with a $proof$, as well as recovering the missing $node$.

Implication Enumeration Given a $theory$, this task requires to enumerate implications of the $theory$, using deductive reasoning.

4.2 Methods

4.2.1 Hypothesis Classification

There are mainly three categories of methods for the hypothesis classification task regarding a multi-task aspect. The first category of methods only conducts the classification task itself; Methods from the second category can predict $correctness$ as well as generate a $proof$. However, the $correctness$ is not necessarily consistent with the predicted $proof$. The third category is similar to the second, except that $correctness$ always follows $proof$.

Until now, methods from the first category directly use transformer-based PLMs (Vaswani et al., 2017), with the target of analyzing and

Method	Generation based	Inference w/ hypothesis	Stepwise	Proof direction	Heuristic search	Verifier	Human-authored realistic proof	Stage
PProver (Saha et al., 2020)	✗	✓	✗	N/A	N/A	✗	✗	1
multiPProver (Saha et al., 2021)	✗	✓	✗	N/A	N/A	✗	✗	1
EntailmentWriter (Dalvi et al., 2021)	✓	✓	✗	N/A	N/A	✗	✓	1
ProofWriter (Tafjord et al., 2021)	✓	✗	✓	→	✗	✗	✗	2
EVR (Liang et al., 2021)	✓	✗	✓	←	✗	✗	✗	2
IBR (Qu et al., 2022)	✗	✓	✓	←	✓	✗	✗	2
IRGR (Ribeiro et al., 2022)	✓	✓	✓	→	✓	✗	✓	2
SI (Creswell et al., 2022)	✓	✗	✓	→	✓	✗	✗	2
FaiRR (Sanyal et al., 2022b)	✓	✗	✓	→	✓	✗	✗	2
MetGen (Hong et al., 2022)	✓	✗	✓	Both	✓	✗	✓	2
SCSearch (Bostrom et al., 2022)	✓	✗	✓	→	✓	✗	✓	2
ADGV (Sprague et al., 2022)	✓	✗	✓	Both	✓	✓	✓	3
NLProofS (Yang et al., 2022a)	✓	✓	✓	→	✓	✓	✓	3
Entailer (Tafjord et al., 2022)	✓	✓	✓	←	✓	✓	✓	3
Teachme (Dalvi et al., 2022)	✓	✓	✓	←	✓	✓	✗	3

Table 2: Methods for Proof Generation task. “Generation based” means whether *proof* is created by generative inference model, otherwise is by utilizing embeddings to classify nodes and edges of *proof*. “Inference w/ hypothesis” means whether *hypothesis* is provided during inference. → and ← denote forward/backward stepwise proof generation. “Heuristic search” with ✗ means exhaustive search. “Human-authored realistic proof” means whether the dataset adopted uses human-authored *proof*, whose contents are consistent with the real world.

benchmarking their performance in different settings (datasets). Specifically, Clark et al. (2020) find that finetuned RoBERTa-large (Liu et al., 2019) can achieve 95%+ accuracy on the test set of D* and ParaRules datasets; Talmor et al. (2020) further demonstrate that LMs can be trained to reliably perform deductive reasoning using both implicit, pre-trained knowledge and explicit natural language statements (*theory*) to make predictions; Han et al. (2022) evaluate finetuned medium-sized language models and few-shot prompting on LLMs on the FOLIO dataset. However, they find that LLM with few-shot prompting only performs slightly better than random results.

The second category methods typically infer PLMs only once, and then utilize the final layer embeddings or generations to obtain *correctness* and *proof*. Specifically, PProver (Saha et al., 2020) and multiPProver (Saha et al., 2021) use the [CLS] token to predict *correctness*, and leverage the final layer embeddings of knowledge items in *theory* to generate *proof*; All-At-Once ProofWriter (Tafjord et al., 2021) and EntailmentWriter (Dalvi et al., 2021) generate *correctness* and linearized *proof* at the same time.

The third category methods create a *proof* first, and then predict *correctness* from the *proof*. §4.2.2 illustrates these methods in detail.

4.2.2 Proof Generation

Current methods for the proof generation task roughly consist of three stages. In each stage, one

key new technique is considered and developed. In stage 1, PLMs are used for forming *proof* in one inference step. In stage 2, modular-based, stepwise frameworks are developed to create *proof* (each module is usually implemented with a single PLM). In stage 3, a verifier is added as a new module to make sure that each reasoning step reflects the belief of PLMs. We will introduce the motivation and typical method for each stage.

Methods for stage 1 typically utilize the last layer embeddings (Saha et al., 2020, 2021) or generations (Tafjord et al., 2021; Dalvi et al., 2021) to create *proof*. Methods utilizing embedding typically (1) obtain an averaged embedding for each knowledge item in *theory*, and (2) pass each embedding to a node classifier, and each embedding pairs to an edge classifier to predict nodes and edges for *proof*. Constraints are usually used to enforce the structure of *proof*. Generation methods directly generate linearized *correctness* and full *proof* given linearized *theory* and *hypothesis*.

The motivations of stage 2 methods are generally concerned with end-to-end methods, which is considered to lack interpretability (Liang et al., 2021; Qu et al., 2022; Sanyal et al., 2022b; Bostrom et al., 2022), suffer from compositional generalization problems (Liang et al., 2021; Creswell et al., 2022), have limited input size (Ribeiro et al., 2022), are not casual (Creswell et al., 2022), and lack constraints on the validity of each inference step (Hong et al., 2022).

311	Methods in stage 2 can be summarized as hav-	4.2.3 Proof with Incomplete Information	363
312	ing two components, an inference module and	ADGV (Sprague et al., 2022) is the only method	364
313	a reasoning controller. The inference module	focusing on this task. It uses both deduction and	365
314	can be a deduction module (Tafjord et al., 2021;	abduction modules, and the reasoning controller	366
315	Ribeiro et al., 2022; Creswell et al., 2022; Sanyal	performs heuristic search. The abduction module	367
316	et al., 2022b; Bostrom et al., 2022), an abduction	is used to recover the missing premise.	368
317	module (Liang et al., 2021; Qu et al., 2022), or		
318	both (Hong et al., 2022; Sprague et al., 2022). The	4.2.4 Implication Enumeration	369
319	deduction module performs deductive reasoning,	Tafjord et al. (2021) is the only paper mentioned	370
320	and reasons forwardly from <i>theory</i> to <i>hypothesis</i>	this task. They compare the performance of “All-	371
321	to construct <i>proof</i> ; the abduction module performs	At-Once” and “Iterative” ProofWriter on this task.	372
322	abductive reasoning, and reasons backwardly from	They find that “All-At-Once” performs worse,	373
323	<i>hypothesis</i> to <i>theory</i> to construct <i>proof</i> . The rea-	mainly because it struggles with problems that are	374
324	soning controller in general performs a search pro-	more complex than training examples.	375
325	cess that each step it searches through the <i>theory</i>		
326	and generated intermediate conclusions space to	4.3 Robustness of PLM as Reasoner	376
327	select (retrieve) premises for the next step infer-	The previously introduced methods only focus on	377
328	ence. The search processes include exhaustive	solving the deductive reasoning tasks, while it is	378
329	search (Tafjord et al., 2021; Liang et al., 2021)	unclear whether PLMs can be used as robust deduc-	379
330	or heuristic search (Qu et al., 2022; Ribeiro et al.,	tive reasoners. To investigate the problem, Gaskell	380
331	2022; Creswell et al., 2022; Sanyal et al., 2022b;	et al. (2022) create a more challenging synthetic	381
332	Bostrom et al., 2022; Hong et al., 2022; Sprague	dataset on hypothesis classification task in terms	382
333	et al., 2022). The reasoning controller usually can	of complexity, and test PLM’s performance on	383
334	also stop the search process if it detects the goal.	it. They find that with large and complex enough	384
335	Motivation of stage 3 methods is similar, basi-	training examples, transformers can perform well	385
336	cally that stage 2 methods lack explicit verifiers	on the dataset. In addition, they find that trans-	386
337	to avoid hallucinating invalid steps (Yang et al.,	formers exhibit some degree of generalization and	387
338	2022a), and to ensure that the inference processes	scale-invariance ability; Richardson and Sabharwal	388
339	reflect PLM’s own beliefs (Tafjord et al., 2022).	(2022) propose an adversarial attack method for	389
340	Methods in stage 3 can be summarized as utiliz-	synthetic datasets on the hypothesis classification	390
341	ing explicit verifier(s) (implemented with a PLM)	task. They find that transformers are often fooled if	391
342	to check the validity of each inference step. One	the query literally appears within the body of a rule,	392
343	way is to add a new module (additional to the infer-	and transformers struggle to correctly bind vari-	393
344	ence module and reasoning controller in stage 2),	ables on either side of a rule; Sanyal et al. (2022a)	394
345	working as a “fact checker” to verify the generated	proposed a synthetic deductive reasoning dataset to	395
346	inference step (Yang et al., 2022a; Tafjord et al.,	evaluate the robustness of language models to min-	396
347	2022); The other one, called round-trip consistency,	imal logical edits in the inputs and different logical	397
348	is only suitable for methods that use both deduc-	equivalence conditions, and find that PLMs are not	398
349	tion and abduction modules, where deduction and	robust to their proposed logical perturbations.	399
350	abduction modules work as the verifier for each		
351	other (Sprague et al., 2022).	5 Inductive Reasoning	400
352	In addition to the general 3 stages, a new aspect	5.1 Existing Task Formulations	401
353	is attended to, which is whether teachable by hu-	Existing tasks for inductive reasoning can be sum-	402
354	mans. Build based on Entailer (Tafjord et al., 2022),	marized as rule verification and rule generation	403
355	TeachMe (Dalvi et al., 2022) shows that user cor-	tasks. Datasets for the tasks are summarized in	404
356	rections can help override erroneous model beliefs,	Table 3. “Generation” tab with \times means it is for	405
357	and that a system can gradually improve by accu-	the rule verification task.	406
358	mulating user corrections. Compared to Entailer, it		
359	adds an interaction module and a dynamic memory	Rule Verification Given a generated <i>rule</i> and	407
360	module to obtain and store human corrections.	<i>facts</i> where the <i>rule</i> is generated from, the task is	408
361	We summarize and analyze the experiment re-	to classify whether the <i>rule</i> can be accepted. The	409
362	sults of proof generation task in §A.7.		

Dataset	Human written	Human labeled	Realistic	Rule provided	Not restricted rule types	Generation	Novel scientific hypotheses	Size
property-norm	✗	✗	✓	✗	✗	✗	✗	23k
DEERLET	✗	✓	✓	✓	✓	✗	✗	846
DEER	✓	✓	✓	✓	✓	✓	✗	1.2k
ARC	✗	✗	-	✗	✗	-	✗	1k
OpenD5	✓	✓	✓	✓	✓	✓	-	675
C-LBD	✓	✗	✓	✓	✓	✓	✓	67k
TOMATO	✓	✓	✓	✓	✓	✓	✓	50

Table 3: Summary of inductive reasoning datasets: property-norm (Misra et al., 2022), DEERLET and DEER (Yang et al., 2022b), ARC (Chollet, 2019), OpenD5 (Zhong et al., 2023), C-LBD (Wang et al., 2023a), and TOMATO (Yang et al., 2023b). “Not restricted rule types” means whether the data is not restricted in a specific topic (e.g., taxonomic).

current evaluation aspects are from requirements of both inductive reasoning and natural language.

Rule Generation Given multiple manually selected *facts* with similar patterns, the task is to induce a *rule* that (1) can entail the *facts*, and (2) is more general than all of the *facts*. Here “more general” means larger information coverage scope. More detailed illustrations can be found in §A.8.

Scientific Hypotheses Generation This task is similar to *Rule Generation* task but is more challenging in that the generated *rule* should not be commonsense knowledge but scientific hypotheses that are even new to humanity.

5.2 Methods

Rule Generation methods almost always have a *Rule Verification* step after the initial generation of rules. To have a clearer overview, we separately introduce the framing or methods of the two tasks.

5.2.1 Rule Verification

Yang et al. (2022b) propose three requirements of rule verification on inductive reasoning from philosophy literature (*rule* and *facts* should not be in conflict; *rule* should reflect reality; *rule* should generalize over *facts*) and one requirement of rule verification from NLP requirement (*rule* should not be trivial or incomplete). They focus on inducing *rule* of many disciplines (e.g., zoology and history) from *facts* as textual observations (e.g. Wikipedia). They implement the verification by LLMs (framing as classification problems).

Another group of works’ (Zhu et al., 2023; Wang et al., 2023b; Qiu and Jiang, 2023) adopted rule verification criteria is compliant with one of the key requirements proposed by Yang et al. (2022b), which is that *rule* and *facts* should not be in conflict. They focus on inducing (executable) *rule* from

synthetic *facts* such as a sequence of number (example *rule*: find the smallest number), arithmetic calculation (example *rule*: “6+4=10”), or changes of 2D grid images (example *rule*: executable code for moving the grids). They verify rules by checking the consistency of the labels of annotated examples (*facts*) and the results of *rules*.

5.2.2 Rule Generation

Yang et al. (2022b) assume that the inductive reasoning task is so difficult that a proper system should contain a rule populator and (multiple) rule verifiers that filter bad rules from different aspects. Accordingly, they propose a framework named chain-of-language-models (CoLM), where one LLM generates *rules* given *facts*, the other four LLMs filter generated rules mainly based on philosophical requirements of inductive reasoning.

Besides the rule generation and filtering process, Zhu et al. (2023) further propose to generate rules based on chain-of-thought prompting, and verify rules based on whether the rules can be used to deduce the annotated answer correctly; Wang et al. (2023b) further propose that under synthetic datasets, executable code can be generated for the textual rules and verify the rules by executing the code and comparing the results with groundtruth annotation; Qiu and Jiang (2023) further propose a third stage of “rule refinement”, and that iteratively repeating the three stages can obtain better rules.

5.2.3 Scientific Hypotheses Generation

Zhong et al. (2023) focuses on proposing hypotheses (from a wide range of disciplines) from a research goal and two comparable corpora. Their method also follows a generate-filter process, where LLMs are used for the filtering stage. Wang et al. (2023a) focus on proposing NLP hypotheses from a seed term and background context. Before hypotheses generation module, they build knowledge graphs to associate academic terms, and retrieve some of the terms as inspirations. Yang et al. (2023b) focuses on proposing social science and business hypotheses only from a pile of raw web corpora. To utilize raw web corpora, they expand generate-filter modules with a background finder module and an inspiration finder module. They also propose three feedback mechanisms named past feedback, present feedback, and future feedback to help the inter-communications between modules to induce more novel, valid, and helpful hypotheses.

Dataset	Human written	Realistic	Multi-step	Theory included	Generation	Size
α NLI	✓	✓	✗	✗	✗	22k
α NLG	✓	✓	✗	✗	✓	76k
AbductionRules	✗	✗	✗	✓	✓	114k
D*-Ab	✗	✗	✓	✓	✓	14k

Table 4: Summary of abductive reasoning datasets: α NLI and α NLG (Bhagavatula et al., 2020), AbductionRules (Young et al., 2022), and D*-Ab (Tafjord et al., 2021). “Realistic” means whether the data is consistent with the real world. “Multi-step” means whether multiple reasoning steps are needed to get the result.

6 Abductive Reasoning

6.1 Existing Task Formulations

Existing tasks for abductive reasoning can be summarized as explanation classification, and explanation generation w/o and w/ theory. Datasets for the tasks are summarized in Table 4. In the table, the “generation” tab and “theory included” tab can be used to determine the task it is used for.

Explanation Classification Given observation O_1 at time t_1 , observation O_2 at time t_2 ($t_2 > t_1$), a plausible hypothesis h^+ and a implausible hypothesis h^- that explain O_1 and O_2 , this task is to select the most plausible hypothesis from h^+ and h^- . O_1 and O_2 each contains a single sentence.

Explanation Generation without Theory

Given observation O_1 at time t_1 , observation O_2 at time t_2 ($t_2 > t_1$), this task is to generate a valid hypothesis h^+ given O_1 and O_2 . O_1 and O_2 each is described in a single sentence.

Explanation Generation with Theory Given a theory C and a possible observation O not provable from C , the task is to generate a new hypothetical fact h such that $C \cup \{h\} \models O$. Here C contains multiple facts and rules, where each fact or rule contains a single sentence. O is in single sentence.

6.2 Methods

6.2.1 Explanation Classification

Methods for this task generally introduce knowledge in various ways to improve performance. Specifically, Mitra et al. (2019) explore ways to incorporate additional unstructured textual knowledge retrieved from a story corpus through prompt; Paul and Frank (2020) encode and incorporate knowledge from COMET’s generation (Bosselut et al., 2019) directly into transformer’s internal attention; Lourie et al. (2021) and Paul and Frank

(2021) incorporate knowledge by multi-task training; Du et al. (2021) incorporate knowledge with an additional pre-training stage using ARL independent story corpora;

In addition to knowledge integration, many different aspects of explanation classification tasks are also investigated. Specifically, Bhagavatula et al. (2020) rewrite the objective using Bayes Rule and formulate a set of probabilistic models that make various independence assumptions on the new objective. They find that the most sophisticated probabilistic model works the best; Zhu et al. (2020) frame this task as a ranking task to also measure the plausibility of hypothesis in addition to discriminating it; Paul and Frank (2021) conduct this task in an unsupervised setting by pretraining on a counterfactual reasoning dataset, which is related to abductive reasoning. Kadikis et al. (2022) propose a method to select suitable PLMs for this task. It is based on the cosine similarity of $embed(O_1, O_2)$ and $embed(h_i)$ for each PLM without finetuning. Zhao et al. (2023) assume that different h are mutually exclusive, and improve performance by incorporating an additional loss item as regularization to enforce an unbalanced probability prediction over different h .

6.2.2 Explanation Generation without Theory

In general, methods for this task either incorporate knowledge or improve the decoding method to be more suitable for this task.

For knowledge integration, Bhagavatula et al. (2020) utilize textual knowledge generated from COMET and investigate two ways of knowledge integration — via texts or via embeddings, and find that the embedding-based method is more effective; Ji et al. (2020) leverage structural knowledge from ConceptNet (Speer et al., 2017) for this task.

For improving decoding method, Qin et al. (2020) are motivated by the fact that the target h^+ to generate happens before O_2 . They accordingly propose an unsupervised decoding algorithm that can incorporate both past and future contexts.

6.2.3 Explanation Generation with Theory

Tafjord et al. (2021) explore the ability of a finetuned T5-11B (Raffel et al., 2020) on $P(h|C, O)$. Their results indicate that finetuned T5-11B can reach a high test accuracy of 93% on D*-Ab.

579
580
581
582
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
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626

7 Challenges of LRNL

Due to the page limit, we list some main challenges in this section, and leave other challenges in §A.9.

Computationally Efficient Reasoner Many tasks in logical reasoning over formal language have very high algorithmic complexity (Muggleton et al., 2012). Thanks to the low computational cost of each deduction step over formal language, such complex tasks could be possible. However, each deduction step in LRNL typically costs one inference of an LLM, which makes tasks with high algorithmic complexity nearly prohibitive.

Robust Reasoner and Reliable Verifier Most methods implement reasoner and verifier with LLMs. It is questionable whether LLMs can robustly reason over any given knowledge. Additionally, the current verifiers only reflect the internal beliefs of LLMs. It is doubtful whether LLMs have obtained the knowledge for verification.

Better Automatic Evaluation Metrics It is generally difficult to automatically evaluate generative reasoning implications, especially with realistic and not synthetic datasets. The difficulty mainly lies in that the same semantic meaning can be expressed with diversified forms, and that different conclusions might be all acceptable (especially in abductive and inductive reasoning). This may lead to biased evaluation when using automatic metrics.

More Impacts on (NLP) Applications As illustrated in §3, overall LRNL can be seen as a new type of neuro-symbolic method, which takes the advantages from both the symbolic and sub-symbolic aspects, and can systematically alleviate many main challenges of both symbolic and sub-symbolic methods. These characteristics make an LRNL system possible (but might still be challenging) to deal with many (NLP) applications such as medical diagnosis and legal NLP tasks, since many medical and legal problems could be seen as pure logical reasoning problems with very large rule bases (e.g., medical knowledge and laws).

Probabilistic Inference In reality, pure deductive reasoning has not always been used. When people include “likely” in their expressions, uncertainty is introduced, which makes the reasoning process probabilistic; In addition, inductive reasoning and abductive reasoning are by default non-monotonic reasoning. This uncertainty aspect has

not been focused in current research. It is probably beneficial to learn from how symbolic reasoning handles uncertainty (Halpern, 2017).

Reasoning with Incomplete Information The current proof generation task requires all necessary premises provided to create a proof tree. Only one work (Sprague et al., 2022) focuses on proof generation with the incomplete information task. However, the task they adopt only overlooks one premise, while in reality more might be missing.

Inductive Reasoning on Web Corpora Currently, the dataset for rule generation tasks in inductive reasoning provides manually selected facts (Yang et al., 2022b). However, to best leverage a system’s ability to handle natural language, it should be able to work on raw web corpora to induce rules, which leads to a more challenging task of inductive reasoning on web corpora.

Abductive Reasoning with (Long) Theory Many tasks such as medical diagnosis conduct abductive reasoning with a long theory (e.g., medical knowledge). However, current abductive reasoning research only covers abductive commonsense reasoning (Bhagavatula et al., 2020) without given theory, or only given short, synthetic, not realistic knowledge as theory (Tafjord et al., 2021).

Interactions between Reasoning Types Multiple reasoning types can be used together for complex tasks. Existing works only utilize deductive reasoning with abductive reasoning to create a proof tree (Hong et al., 2022; Sprague et al., 2022). However, many other collaborations are possible, such as using inductive reasoning to collect a (large) rule base, which is to be used as the theory base for deductive reasoning.

8 Conclusion

In this paper, we summarize the three previously separately investigated logical reasoning types together, referred as logical reasoning from the perspectives of deductive, inductive, and abductive reasoning over natural language as knowledge representation (LRNL), and provide an in-depth and up-to-date survey of LRNL. Specifically, we have introduced the philosophical foundations, advantages of LRNL, benchmarks and methods, challenges of LRNL, possible future directions, and the relation of LRNL to related NLP fields (§A.1).

627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673

9 Limitations

In consideration of space constraints, this paper focuses more on (1) providing a high-level overview and prospect of the LRNL field (e.g., advantages and challenges of the field), and (2) delineating the broader evolutionary trajectories of pertinent methodologies. It might not include all the details of the surveyed papers.

10 Ethics Statement

This article follows the ACL Code of Ethics. To our knowledge, there are no foreseeable potential risks to use the datasets and methods in this paper.

References

- A. Aamodt and E. Plaza. 1994. Case-based reasoning:foundational issues,methodological variations,and system approaches. *AI communications*.
- Qiming Bao, Alex Yuxuan Peng, Tim Hartill, Neset Tan, Zhenyun Deng, Michael Witbrock, and Jiamou Liu. 2022. Multi-step deductive reasoning over natural language: An empirical study on out-of-distribution generalisation. The 2nd International Joint Conference on Learning and Reasoning and 16th International Workshop on Neural-Symbolic Learning and Reasoning (IJCLR-NeSy 2022).
- 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 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. [COMET: commonsense transformers for automatic knowledge graph construction](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4762–4779. Association for Computational Linguistics.
- Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, and Greg Durrett. 2022. [Natural language deduction](#)

[through search over statement compositions](#). *CoRR*, abs/2201.06028.

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 632–642. The Association for Computational Linguistics.
- Erik Cambria, Qian Liu, Sergio Decherchi, Frank Xing, and Kenneth Kwok. 2022. [Senticnet 7: A commonsense-based neurosymbolic AI framework for explainable sentiment analysis](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 3829–3839. European Language Resources Association.
- Erik Cambria, Lorenzo Malandri, Fabio Mercurio, Mario Mezzanzanica, and Navid Nobani. 2023. [A survey on XAI and natural language explanations](#). *Inf. Process. Manag.*, 60(1):103111.
- François Chollet. 2019. [On the measure of intelligence](#). *CoRR*, abs/1911.01547.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. [Transformers as soft reasoners over language](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3882–3890. ijcai.org.
- Luca Console and Lorenza Saitta. 2000. On the relations between abductive and inductive explanation. *Abduction and Induction: Essays on their Relation and Integration*, pages 133–151.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. [Selection-inference: Exploiting large language models for interpretable logical reasoning](#). *CoRR*, abs/2205.09712.
- Andrew Cropper, Sebastijan Dumancic, Richard Evans, and Stephen H. Muggleton. 2022. [Inductive logic programming at 30](#). *Mach. Learn.*, 111(1):147–172.
- Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. [Explaining answers with entailment trees](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 7358–7370. Association for Computational Linguistics.
- Bhavana Dalvi, Oyvind Tafjord, and Peter Clark. 2022. [Towards teachable reasoning systems](#). *CoRR*, abs/2204.13074.
- Rajarshi Das, Ameya Godbole, Ankita Naik, Elliot Tower, Manzil Zaheer, Hannaneh Hajishirzi, Robin Jia, and Andrew McCallum. 2022. [Knowledge base question answering by case-based reasoning over](#)

1115	Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope.	Nathan Young, Qiming Bao, Joshua Bensemann, and	1171
1116	2023a. Learning to generate novel scientific directions with contextualized literature-based discovery.	Michael Witbrock. 2022. Abductionrules: Training transformers to explain unexpected inputs. In <i>Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 218–227. Association for Computational Linguistics.	1172
1117	<i>CoRR</i> , abs/2305.14259.		1173
1118			1174
1119	Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah D. Goodman. 2023b. Hypothesis search: Inductive reasoning with language models. <i>CoRR</i> , abs/2309.05660.	Fei Yu, Hongbo Zhang, and Benyou Wang. 2023. Nature language reasoning, a survey. <i>arXiv preprint arXiv:2303.14725</i> .	1177
1120			1178
1121			1179
1122			
1123	Wenguan Wang and Yi Yang. 2022. Towards data- and knowledge-driven artificial intelligence: A survey on neuro-symbolic computing. <i>arXiv preprint arXiv:2210.15889</i> .	Wenting Zhao, Justin T Chiu, Claire Cardie, and Alexander M Rush. 2023. Abductive commonsense reasoning exploiting mutually exclusive explanations. <i>arXiv preprint arXiv:2305.14618</i> .	1180
1124			1181
1125			1182
1126			1183
1127	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. <i>CoRR</i> , abs/2201.11903.	Ruiqi Zhong, Peter Zhang, Steve Li, Jinwoo Ahn, Dan Klein, and Jacob Steinhardt. 2023. Goal driven discovery of distributional differences via language descriptions. <i>CoRR</i> , abs/2302.14233.	1184
1128			1185
1129			1186
1130			1187
1131	Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. 2023. Are large language models really good logical reasoners? a comprehensive evaluation from deductive, inductive and abductive views. <i>arXiv preprint arXiv:2306.09841</i> .	Yunchang Zhu, Liang Pang, Yanyan Lan, and Xueqi Cheng. 2020. L2r²: Leveraging ranking for abductive reasoning. In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 1961–1964. ACM.	1188
1132			1189
1133			1190
1134			1191
1135			1192
1136	Kaiyu Yang, Jia Deng, and Danqi Chen. 2022a. Generating natural language proofs with verifier-guided search. <i>CoRR</i> , abs/2205.12443.		1193
1137			1194
1138			
1139	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2369–2380. Association for Computational Linguistics.	Zhaocheng Zhu, Yuan Xue, Xinyun Chen, Denny Zhou, Jian Tang, Dale Schuurmans, and Hanjun Dai. 2023. Large language models can learn rules. <i>arXiv preprint arXiv:2310.07064</i> .	1195
1140			1196
1141			1197
1142			1198
1143			
1144			
1145			
1146			
1147			
1148	Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022b. Language models as inductive reasoners. <i>CoRR</i> , abs/2212.10923.	A Appendix	1199
1149			
1150			
1151			
1152	Zonglin Yang, Xinya Du, Erik Cambria, and Claire Cardie. 2023a. End-to-end case-based reasoning for commonsense knowledge base completion. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3509–3522, Dubrovnik, Croatia. Association for Computational Linguistics.	A.1 Relation to Related (NLP) Fields	1200
1153			
1154			
1155			
1156			
1157			
1158			
1159	Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. 2023b. Large language models for automated open-domain scientific hypotheses discovery. <i>CoRR</i> , abs/2309.02726.	In this section, we first introduce related NLP fields to general logical reasoning, then introduce fields that are only related to deductive reasoning, inductive reasoning, or abductive reasoning. We hope that this section could be helpful to form a clear shape of LRNL in NLP.	1201
1160			1202
1161			1203
1162			1204
1163			1205
1164			1206
1165			
1166			
1167			
1168			
1169			
1170			
		A.1.1 Logical Reasoning	1207
		There are some previous works involve the term “logical reasoning”, but do not provide a specification on which sub-type of logical reasoning they involve. In many cases these works are more close to “natural language inference”, which adopts datasets where the data involve a mixture of multiple sub-types of logical reasoning, making it hard to analyze from each sub-type. Therefore we do not include these works in this survey.	1208
			1209
			1210
			1211
			1212
			1213
			1214
			1215
			1216
		Neuro-Symbolic Computing Neural-symbolic computing (NeSy) is a hybrid of symbolism and connectionism to exploit advantages from both sides (Wang and Yang, 2022; Cambria et al., 2022). The knowledge representation of its symbolic part	1217
			1218
			1219
			1220
			1221

1222	basically is a knowledge graph or propositional	size (Ribeiro et al., 2022), and contains unrelated	1272
1223	logic or first-order logic (Wang and Yang, 2022).	or incorrect steps (Hong et al., 2022; Tafjord et al.,	1273
1224	LRNL could be seen as a new type of NeSy in addition	2022).	1274
1225	to the existing 6 types summarized by Kautz	Overall, it could be interesting to use COT-	1275
1226	(2022), as its goal and design of methodology are	related methods specifically for deductive, inductive,	1276
1227	typically symbolic (logical reasoning with knowl-	edge, or abductive reasoning (as opposed to modular-	1277
1228	edge bases), while avoiding any symbolic represen-	based methods), and it is a less-explored research	1278
1229	tation, using (currently pure) neural methods.	direction.	1279
1230	Natural Language Inference Natural language	A.1.2 Deductive Reasoning	1280
1231	inference (NLI) is generally considered as the	Multi-hop Reasoning Compared to proof gen-	1281
1232	semantic concepts of entailment and contradic-	eration, many multi-hop reasoning tasks (Yang	1282
1233	tion (Bowman et al., 2015). Here logical reasoning	et al., 2018; Jiang et al., 2020; Min et al., 2019;	1283
1234	tasks can be viewed as special types of NLI focus-	Sinha et al., 2019) are much simpler, often being	1284
1235	ing on particular reasoning aspects.	single-branched (Qu et al., 2022), consisting of	1285
1236	Question Answering The form of LRNL looks	only 2-3 supporting facts, and are more coarse-	1286
1237	similar to question answering (QA), however, QA	grained, involving large chunks of texts such as	1287
1238	is conducting one-step logical reasoning only when	passages instead of simple, short sentences (Yang	1288
1239	the context provides enough information to answer	et al., 2022a).	1289
1240	the question (deductive reasoning), or the answer	Nevertheless, some multi-hop reasoning datasets	1290
1241	is a generalization of an argument in context or	can be considered as conducting deductive reason-	1291
1242	question (inductive reasoning), or the answer is	ing. For instance, for each data in CLUTRR (Sinha	1292
1243	to provide explanations to the question (abductive	et al., 2019) dataset, a set of facts that can make	1293
1244	reasoning).	conclusive support to the target kinship relation is	1294
1245	Commonsense Reasoning Commonsense reason-	included in background information as input for	1295
1246	ing (CR) and logical reasoning (LR) are similar	each target relation, hence from the philosophical	1296
1247	in that they both involve “knowledge” and “reason-	definition (Salmon, 1989), it requires to perform	1297
1248	ing”. Compared to LR, CR focuses more on the	deductive reasoning.	1298
1249	“knowledge” aspect. Some typical tasks include	Mathematical Reasoning In many mathemati-	1299
1250	whether a system has commonsense knowl-	cal reasoning tasks such as math word problem	1300
1251	edge (Bosselet et al., 2019; Yang et al., 2020),	solving (Koncel-Kedziorski et al., 2015) and geom-	1301
1252	and whether a system’s answer is commonsense-	etry problem solving (Seo et al., 2015), the conclu-	1302
1253	knowledge-aware (Bisk et al., 2020); LR focuses	sion can be conclusively entailed by the premise.	1303
1254	more on the “reasoning” aspect, e.g., whether a	Therefore these tasks belong to deductive reason-	1304
1255	system’s i/o behaviors follow reasoning require-	ing. We do not review math-related papers because	1305
1256	ments (Clark et al., 2020).	we want to focus solely on the challenge of de-	1306
1257	Chain of Thoughts Chain of	ductive reasoning while mathematical reasoning	1307
1258	thoughts (COT) (Wei et al., 2022) is a prompting	involves numbers in the text, which introduces ad-	1308
1259	technique that can elicit the step-by-step reasoning	ditional challenges.	1309
1260	ability of LLMs without finetuning.	A.1.3 Inductive Reasoning	1310
1261	COT can potentially be used for each of the three	Information Extraction Information Extrac-	1311
1262	sub-reasoning types of logical reasoning. In fact,	tion (IE) is a task of extracting pre-specified types	1312
1263	for a given (commonsense reasoning) question,	of facts from written texts or speech transcripts, and	1313
1264	some reasoning steps of COT could be deductive,	converting them into structured representations (Ji,	1314
1265	and others can be inductive or abductive. Since the	2018). The rule generation task here also extracts	1315
1266	purpose of this paper is to provide a finer analysis	rules from facts represented in written texts. The	1316
1267	on logical reasoning, we do not intentionally cover	difference is that IE pursues extracting the exact	1317
1268	prompting techniques such as COT.	information from existing texts, while inductive	1318
1269	It is also argued by several modular-based de-	reasoning aspires to induce more general rules from	1319
1270	ductive reasoning methods that COT’s reasoning is	existing texts, where the information in rules goes	1320
1271	not casual (Creswell et al., 2022), limited by input	beyond what is exactly stated in the texts.	1321

Case-based Reasoning Case-based Reasoning (CBR) is a classic AI subject, whose methods share a general methodology of four steps: retrieve, reuse, revise, and retain (Aamodt and Plaza, 1994). Recently there has been research works devoting to bridging the research of CBR and NLP, by using NLP techniques for CBR challenges (Yang et al., 2023a) and improving NLP tasks with CBR methodologies (Das et al., 2021, 2022; Yang et al., 2023a; Thai et al., 2023). CBR could be seen as a type of analogical reasoning (Kolodner, 1997), and analogical reasoning belongs to inductive reasoning (Salmon, 1989). However, CBR is a different inductive reasoning type than the “generalization” process (from facts to rules) described in Flach and Kakas (2000), but more on the general description on inductive reasoning (Salmon, 1989) that premises cannot conclusively provide support to the conclusion.

A.1.4 Abductive Reasoning

Casual Reasoning In logic research, causal reasoning aims at an epistemological problem of establishing precise causal relationships between causes and effects. It is generally considered a form of inductive reasoning (Goertzel et al., 2011), since inductive reasoning is to derive rules that lead from one to another. When the focus is to derive possible causes from effects, the problem belongs to abductive reasoning (Goertzel et al., 2011).

A.2 More Details About the Definition and Categorization of Logical Reasoning

There are many subjects related to logical reasoning, including philosophy, logic, and AI. Among them, the definition and categorization aspects of logical reasoning are handled by philosophy research. However, debate exists in philosophy research on the categorization of logical reasoning.

One group believes that every argument can be classified as either deduction argument, inductive argument, or fallacy (Salmon, 1989). Without considering fallacy, given that an argument consists of premises and a conclusion, when the premises can conclusively provide support to the conclusion (which means that if the premises of the argument were all true, it would be impossible for the conclusion of the argument to be false), this argument is a deductive argument. Conversely, when the premises can not conclusively provide support to the conclusion, the argument is inductive.

The other group has the same definition of de-

ductive reasoning, but they believe that further categorization of non-deductive reasoning is necessary. Without considering fallacy, they believe in a trichotomy of deductive, inductive, and abductive reasoning (Peirce, 1974). However, even for the second group, the definition and difference between inductive and abductive reasoning are also controversy (Flach and Kakas, 2000).

Nevertheless, Console and Saitta (2000) argue that from the utility perspective of AI, a distinction between inductive and abductive reasoning is possible: both inductive and abductive reasoning provide explanations about the world but their explanations differ in the degree of generality. For instance, an inductive hypothesis allows the validity of properties, observed on a set of individuals, to be generalized to other individuals not in the observations, whereas an abductive one allows unobserved properties to be applied to observed individuals. More details about the difference and an example can be found in §A.2.

Considering that inductive and abductive reasoning can be distinctive enough when formulated in NLP, in this paper, we adopt the second group, particularly Console and Saitta (2000)’s view of definition and categorization of logical reasoning.

Specifically, the difference between inductive and abductive reasoning is that, both inductive and abductive reasoning provide explanations about the world but their explanations differ in the degree of generality.

For instance, an inductive hypothesis allows the validity of properties, observed on a set of individuals, to be generalized to other individuals not in the observations, whereas an abductive one allows unobserved properties to be applied to observed individuals.

The distinction between inductive and abductive hypotheses strictly parallels the dichotomy *extension* vs. *intension*, or *generality* vs. *informativeness*. In other words, an inductive hypothesis extends or generalizes to unobserved individuals, while an abductive one provides more specific information (e.g., unobserved properties) about existing specific individuals.

For example, if a white ball is found in a bag, inductive reasoning might lead to the conclusion that “all balls in this bag are white”, while abductive reasoning might lead to the conclusion that “someone put the white ball into this bag”.

In this example, the inductive hypothesis generalizes the property of the existing individual (a

1424	found white ball) to unobserved individuals (other	1475
1425	not-seen balls in the bag), while the abductive hy-	1476
1426	pothesis provides more specific information about	1477
1427	the current individual (who brought this ball to the	1478
1428	bag).	1479
1429	To summarize in simple words, in common sit-	1480
1430	uations, pure inductive reasoning is to only pro-	1481
1431	vide (usually sample to population) generalizations,	1482
1432	while pure abductive reasoning is to only provide	1483
1433	specific explanations.	1484
1434	Overall, even in the philosophical literature	1485
1435	(which takes charge of the research on the defini-	
1436	tion of logical reasoning), a clear definition for all	1486
1437	three types of logical reasoning is rare, but more on	
1438	the description of the difference between types of	1487
1439	logical reasoning (since a clear definition is still un-	1488
1440	der debate). The difference can be illustrated does	1489
1441	not mean a precise definition can be given. Never-	1490
1442	theless, considering the above-discussed philosoph-	1491
1443	ical literature, we try our best to give a definition	1492
1444	below for a more straightforward understanding:	1493
1445	Given an argument consisting of premises and	1494
1446	a conclusion, we define the sub-type of logical	1495
1447	reasoning it involves below:	1496
1448	<i>Definition for deductive reasoning:</i> the premises	
1449	can conclusively provide support for the conclu-	1497
1450	sion, i.e. if the premises are all true, it would be	1498
1451	impossible for the conclusion to be false.	
1452	<i>Definition for inductive reasoning:</i> the premises	1499
1453	cannot conclusively provide support for the conclu-	1500
1454	sion, since the conclusion generalizes existing	1501
1455	information in premises to new knowledge, which	1502
1456	has a wider applicable scope than those in premises.	1503
1457	<i>Definition for abductive reasoning:</i> the premises	1504
1458	cannot conclusively provide support for the conclu-	1505
1459	sion, since the conclusion contains more specific	1506
1460	information over the premises (most commonly	1507
1461	used as generating most probable explanations).	1508
1462	Please note that according to Console and Saitta	1509
1463	(2000), inductive reasoning and abductive reason-	1510
1464	ing are not exclusive to each other, i.e., inductive	1511
1465	reasoning and abductive reasoning overlap with	1512
1466	each other.	1513
1467		1514
1468	A.3 Related Surveys on Reasoning	1515
1469	Huang and Chang (2022) ; Qiao et al. (2022) mainly	1516
1470	reviews the prompting techniques for LLMs, but	1517
1471	do not focus on papers that specialized on logical	1518
1472	reasoning (the coverage of the two fields are quite	1519
1473	different).	1520
1474	Yu et al. (2023) is a concurrent work of ours and	1521
	reviews papers related to reasoning. However, it	1522
		1523
		1524
	does not focus on logical reasoning, particularly	
	the three sub-types of logical reasoning. The advan-	
	tage of our survey is that we provide a finer analysis	
	of logical reasoning (including a more detailed def-	
	inition and categorization of logical reasoning from	
	philosophy literature, comparison with the classic	
	AI paradigm on logical reasoning, and organizing	
	the survey based on the three sub-types of logical	
	reasoning).	
	Xu et al. (2023) provides a comprehensive eval-	
	uation of the logical reasoning ability of LLMs.	
	A.4 Other Related Logical Reasoning Papers	
	A.4.1 Rule Verification	
	Misra et al. (2022) analyze language model’s abil-	
	ity to generalize novel property knowledge (has	
	sesamoid bones) from concept(s) (robins) to oth-	
	ers (sparrows, canaries). As illustrated in §A.8,	
	they analyze the language models’ ability to clas-	
	sify a new fact (but not a rule) as correct or not,	
	given facts. It could be seen that the correctness	
	of a rule is implicitly predicted by testing multiple	
	facts entailed by the rule.	
	A.5 Research Trend in the Three Sub-Types	
	of Logical Reasoning	
	Out of the three reasoning types, deductive reason-	
	ing has drawn the most research attention, and has	
	the most abundant of works, especially in 2022.	
	Abductive reasoning has drawn much attention in	
	2020 and 2021 but has few works in 2022 and 2023.	
	Inductive reasoning is only proposed at the end of	
	2022, having the least number of works. However,	
	inductive reasoning has attracted much attention	
	since the second half year of 2023.	
	Two main reasons for the abundance of works	
	in the deductive reasoning domain could be that	
	(1) more challenging benchmarks have been con-	
	structed during the last few years, and (2) deductive	
	reasoning could be one of the most commonly used	
	reasoning types in common life. We think the main	
	reason for the little attention drawn to abductive	
	reasoning in recent years is that the benchmarks for	
	abductive reasoning are relatively old and less chal-	
	lenging for LLMs. Inductive reasoning could be a	
	promising research topic since there have been few	
	works in the domain, and it involves very challeng-	
	ing tasks such as proposing new scientific findings.	
	In general, there has been no framework which	
	is proposed to address all three reasoning domains.	
	However, LLMs generally can exhibit all three rea-	
	soning abilities to some extent. It would be inter-	

esting for future works to analyze the effect of the pretraining method and scale of PLM on the three reasoning abilities.

A.6 Relation Between LRNL and NeSy

A large proportion of recent papers on deductive reasoning and abductive reasoning leverage a natural language-based knowledge base, and reason over retrieved knowledge from the knowledge base to reach a certain goal (Tafjord et al., 2021; Liang et al., 2021; Qu et al., 2022; Ribeiro et al., 2022; Creswell et al., 2022; Sanyal et al., 2022b; Hong et al., 2022; Bostrom et al., 2022; Yang et al., 2022a; Tafjord et al., 2022; Dalvi et al., 2022). This pattern is very similar to the methodology design of NeSy, which is to retrieve symbolic knowledge and reason over the retrieved symbolic knowledge. The main difference is that LRNL adopts natural language as knowledge representation but not symbolic knowledge. Because of the similarity in the methodology design, we consider that LRNL could be seen as a type of NeSy but without many disadvantages of symbolic representation such as symbolic knowledge acquisition and scalability.

In addition, due to the high similarity in the methodology design to NeSy, LRNL also shares some advantages with NeSy such as explainability. The reason is that the iterative retrieving and reasoning will make the decision-making process more interpretable on the intermediate reasoning steps, and which knowledge is used for each reasoning step.

A.7 Experiments Summarization

In this section, we summarize the experiment results of an important and literature-abundant task.

Until now there has been only one or two papers working on inductive reasoning. Methods for abductive reasoning generally leverage different resources (such as multi-task, additional knowledge resources, and ancillary loss) and lack a progressive relationship between each other, therefore are less comparable. Currently, the *ProofGeneration* task in deductive reasoning is the most literature-abundant, and methods for this task have progressive relationships with each other. Therefore here we mainly summarize results and analyze for the *ProofGeneration* task.

Table 5 shows the summarized experiment results. We select the most widely used tasks to display their performance. Among the task, the setting of ParaRules is trained on D3 (D* dataset

with depth 3) and tested on the ParaRules test set; the setting of Birds-Electricity is trained on D5 (D* dataset with depth 5) and tested on bird-electricity set; setting for EntailmentBank is the task 3 which uses full corpus as input (so that many distractors exist in input); setting for OBQA and QuaRTz are zero-shot setting while model pre-trained on another dataset (EntailmentBank).

Among the methods, Creswell et al. (2022) and Bostrom et al. (2022) design unique metrics using EntailmentBank dataset, and Sprague et al. (2022) focus on a unique task (proof generation task with incomplete information), therefore we do not list their experiments results in the table.

Overall methods for proof generation tasks tend to use different datasets for evaluation, making them less comparable.

A.8 Meaning of “More General” Required by Inductive Reasoning

This section is collected from Yang et al. (2022b)’s appendix, to help illustrate inductive reasoning.

Given an argument consisting of a premise and a conclusion, if the conclusion involves new information that is not covered by the premise and can not be conclusively entailed by the premise, the argument is an inductive argument (Salmon, 1989).

When the conclusion has a larger scope of information coverage than the premise, and can entail the premise, it can be said that the conclusion is “more general” to the premise (Yang et al., 2022b). In this case, we termed the premise as a “fact”, and the conclusion as a “rule”; When the conclusion contains new pieces of information and cannot entail the premise, as defined by Salmon (1989), the argument is still an inductive argument. But in this case, we termed the premise as a “fact”, and the conclusion as another “fact”.

For instance, if facts that are about cats and dogs are good accompaniment of humans, then some examples of a “more general” rule can be (1) mammals are good accompaniment of humans, or (2) domesticated animals are good accompaniment of humans, or (3) animals with four legs are good accompaniment of human.

In these examples, the rules cover a larger scope than the facts (e.g., mammals compared to cats; domesticated animals compared to cats), and therefore the rules are “more general” than the facts.

“More general” means not only about finding higher taxonomic rank, but can be in unlimited forms. For instance, if the fact is about the Sun

Methods	ParaRules	Birds-Electricity	EntailmentBank (Task 3)							OBQA	QuaRTz
	Full Accuracy (FA)	Full Accuracy (FA)	Leaves F1	Leaves All-Cor.	Steps F1	Steps All-Cor.	Intermediates F1	Intermediates All-Cor.	Overall All-Correct	Accuracy	Accuracy
PProver	95.1	80.5	-	-	-	-	-	-	-	-	-
multiPProver	94.5	81.8	-	-	-	-	-	-	-	-	-
EntailmentWriter	-	-	39.7	3.8	7.8	2.9	36.4	13.2	2.9	-	-
ProofWriter	98.5	97.0	-	-	-	-	-	-	-	-	-
EVR	-	63.1	-	-	-	-	-	-	-	-	-
IBR	95.7	93.5	-	-	-	-	-	-	-	-	-
IRGR	-	-	45.6	12.1	16.3	11.8	38.8	36.5	11.8	-	-
Selection-Inference	-	-	-	-	-	-	-	-	-	-	-
FaiRR	98.6	-	-	-	-	-	-	-	-	-	-
MetGen	-	-	34.8	8.7	9.8	8.6	36.7	20.4	8.6	-	-
SCSearch	-	-	-	-	-	-	-	-	-	-	-
ADGV	-	-	-	-	-	-	-	-	-	-	-
NLProofS	-	-	43.2	8.2	11.2	6.9	42.9	17.3	6.9	-	-
Entailer	-	-	-	-	-	-	-	-	-	76.8	74.3
Teachme	-	-	-	-	-	-	-	-	-	77.0	75.9

Table 5: Proof Generation Task Results.

rises and falls every day, then some examples of a “more general” rule can be (1) the Earth is the king of the universe or (2) the Earth is rotating itself.

Both rule examples are “more general” than the given fact, since the rule can entail not only the given fact, but also other not mentioned facts such as the observable movements of the other stars in the Milky Way.

A.9 Other Challenges and Possible Future Directions

Robust Deductive Reasoner Symbolic deductive reasoners are not restricted to train data distributions, while neural deductive reasoners are restricted to their training data (Gontier et al., 2020; Richardson and Sabharwal, 2022); In addition, neural deductive reasoners are also vulnerable to adversarial attacks (Gaskell et al., 2022), while symbolic reasoners are robust to the attacks. The lack of robustness can lead to restricted application domains and incorrect deductive inferences.

Reliable Rule Generation Currently, the rule generation method in inductive reasoning relies on out-of-box LLMs, since a finetuned rule generation model could be restricted in a domain. The annotation of an inductive reasoning dataset should only be done by experts and is very time consuming (Yang et al., 2022b). Given the two restrictions, how to improve the quality of generated rules given related facts could be a challenging open problem.

Reliable Explanation Generation Abduction is a form of non-monotonic reasoning (Paul, 1993), and potentially has a large search space of conclusions given premises. Therefore, how to generate more (all) reasonable explanations can be challenging (Bhagavatula et al., 2020).

Building Larger Benchmarks For complicated reasoning tasks especially in realistic and natural language settings, usually experts are needed

for annotation, and the process is very time-consuming (Dalvi et al., 2021; Sprague et al., 2022; Yang et al., 2022b). Therefore it can be challenging to construct significantly larger benchmarks.

Understanding the Internal Mechanism of PLMs for Reasoning Until now research works only focused on investigating whether the input/output behaviors of PLMs can be used to simulate a reasoner (Clark et al., 2020) or complete reasoning tasks. However, it is still a challenging open research question to understand the internal mechanism of PLMs for reasoning.