

Autoformalizing Natural Language to First-Order Logic: A Case Study in Logical Fallacy Detection

Anonymous ACL submission

Abstract

Translating natural language into formal language such as First-Order Logic (FOL) is a foundational challenge in NLP with wide-ranging applications in automated reasoning, misinformation tracking, and knowledge validation. In this paper, we introduce Natural Language to First-Order Logic (NL2FOL), a framework to autoformalize natural language to FOL step-by-step using Large Language Models (LLMs). Our approach addresses key challenges in this translation process, including the integration of implicit background knowledge. By leveraging structured representations generated by NL2FOL, we use Satisfiability Modulo Theory (SMT) solvers to reason about the logical validity of natural language statements. We present logical fallacy detection as a case study to evaluate the efficacy of NL2FOL. Being neurosymbolic, our approach also provides interpretable insights into the reasoning process and demonstrates robustness without requiring model fine-tuning or labeled training data. Our framework achieves strong performance on multiple datasets – on the LOGIC dataset, NL2FOL achieves an F1-score of 78%, while generalizing effectively to the LOGIC-CLIMATE dataset with an F1-score of 80%.

1 Introduction

In recent years, Large Language Models (LLMs) have shown impressive advancements in understanding and generating natural language (Brown et al., 2020). Despite this progress, their ability to tackle complex reasoning tasks remains limited (Bubeck et al., 2023; Wei et al., 2022). These challenges are especially prevalent in multistep logical deductions, abstract reasoning, and knowledge integration in various domains (Dalvi et al., 2021; Chen et al., 2024). Addressing these limitations and improving the reasoning capabilities of LLMs has become a critical focus in AI research (Haluptzok et al., 2022; Gendron et al., 2024).

In contrast, formal reasoning tools such as Satisfiability Modulo Theory (SMT) solvers excel in reasoning, providing rigorous, provable guarantees by leveraging symbolic representations and logical calculus (Barrett et al., 2009; De Moura and Bjørner, 2008). However, a key limitation of formal solvers is their reliance on structured logical input, such as First Order Logic (FOL), which must accurately capture the semantics and context of natural language statements (Beltagy et al., 2016). This presents the challenge of translating unstructured natural language into a structured form required for formal reasoning while preserving essential context and meaning.

This also brings a unique opportunity: if we can reliably translate natural language into structured logical forms, we can harness the power of formal solvers to reason systematically over natural language statements. However, achieving this translation is nontrivial, as it involves accurately capturing natural language semantics (Beltagy et al., 2016). Moreover, translating to a formal logical form may cause implicit and external context to be lost, which must be reintroduced to ensure logical accuracy.

To address these challenges, we present NL2FOL, a novel framework that bridges the gap between natural language and formal reasoning systems. NL2FOL employs a structured, step-by-step pipeline to translate natural language inputs into first-order logic (FOL) representations, leveraging large language models (LLMs) at each step for enhanced precision and adaptability. A distinguishing feature of NL2FOL is its seamless integration of background knowledge into the generated logical forms, overcoming a major limitation of traditional formal logic frameworks - the inability to capture implicit information embedded in natural language.

In this paper, we demonstrate the effectiveness of NL2FOL through a case study on logical fallacy

Fallacy Name	Example	Logical Form
Faulty Generalization	Sometimes flu vaccines don't work; therefore vaccines are useless.	$(\exists x \in \text{FluVaccines}(\text{DoesntWork}(x)) \wedge (\text{FluVaccines} \subseteq \text{Vaccines})) \Rightarrow (\forall y \in \text{Vaccines}(\text{DoesntWork}(y)))$
False Causality	Every time I wash my car, it rains. Me washing my car has a definite effect on the weather.	$\text{occuredAfter}(\text{washingCar}, \text{rain}) \Rightarrow \text{caused}(\text{washingCar}, \text{rain})$
Ad Populum	Everyone should like coffee: 95% of teachers do!	$(\text{like}(\text{coffee}, 95\% \text{Teachers})) \Rightarrow (\text{like}(\text{coffee}, \text{everyone}))$
False Dilemma	I don't want to give up my car, so I don't think I can support fighting climate change.	$\forall(a)(\text{giveUpCar}(a) \vee \text{dontSupportFightingClimateChange}(a))$

Table 1: Sample logical fallacies from Jin et al. (2022) along with examples and their logical forms. For each type of fallacy, we show one possible logical form.

detection, showcasing its ability to identify and explain faulty reasoning in natural language arguments. Detecting logical fallacies is particularly challenging as they often rely on reasoning patterns that appear plausible yet are fundamentally flawed (Jin et al., 2022). To address this, NL2FOL translates logical fallacies from natural language into FOL representations, enabling formal solvers to verify logical validity. These solvers generate counterexamples and explanations, which are interpreted back into natural language to enhance human comprehensibility. By incorporating intermediate natural language outputs, our pipeline improves interpretability, transparency, and debuggability (Bai et al., 2020).

We show that our framework achieves strong performance on the logical fallacy detection benchmarks LOGIC and LOGICCLIMATE (Jin et al., 2022), with F1 scores of 78% and 80%, respectively - outperforming existing models by 22% on the challenge set, LOGICCLIMATE. These results highlight NL2FOL as a generalizable and interpretable tool for reasoning tasks that demand the precision of formal reasoning systems.

By analyzing the strengths and weaknesses of LLMs at each step of the NL2FOL pipeline, we further identify opportunities for improving logical reasoning capabilities. Even though LLMs prove to be effective in parsing and generating logical representations for structured inputs, they often struggle with ambiguities in natural language and incorporating nuanced contextual knowledge. The ability to integrate symbolic solvers with language models positions NL2FOL as a powerful neurosymbolic approach, bridging the gap between formal reasoning and natural language understanding.

2 Related Work

Logical fallacy detection. Existing work on classifying logical fallacies includes argument sufficiency classification (Stab and Gurevych, 2017), ad hominem fallacies from Reddit posts (Habernal et al., 2018b) and dialogues (Habernal et al., 2018a), rule parsers (Nakpih and Santini, 2020), structure-aware Transformers (Jin et al., 2022), multitask instruction based prompting (Alhindi et al., 2022), and instance-based reasoning (Sourati et al., 2022). To our knowledge, our work is the first on few-shot classification of logical fallacies in a step-by-step, explainable manner. By ensuring that the reasoning process is transparent, we allow users to understand and verify the system decision.

Natural language to formal logic. While early work on mapping text to formal logic relied heavily on grammar-based approaches (Purdy, 1991; Angeli and Manning, 2014; MacCartney and Manning, 2014), recent advances in deep learning and foundation models have enabled new data-driven techniques for translating natural language to linear temporal logic (Cosler et al., 2023; Fuggitti and Chakraborti, 2023; Liu et al., 2022) and first-order logic (Singh et al., 2020; Yang et al., 2024; Hahn et al., 2022). Neural models for parsing natural language to first-order logic (Singh et al., 2020; Yang et al., 2024) and neuro-symbolic approach combining language models with first-order logic provers (Olausson et al., 2023) have since been explored. However, these approaches still face challenges in accurately capturing implicit information or transforming complex ambiguous sentences into logical form, mainly attributed to linguistic ambiguity.

Aly et al. (2023) integrated LLMs with logical inference for fact verification, and while our method

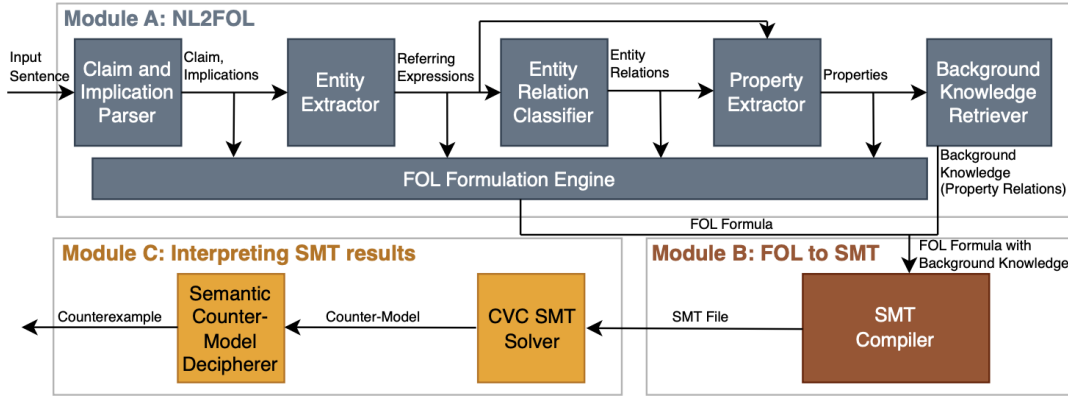


Figure 1: Overview of the proposed framework used for logical fallacy detection. *Module A* converts natural language input to a first-order logic formula merged with contextual relationships, *Module B* compiles the negation of a given logical formula to an SMT file with well-defined sorts for variables and predicates, and *Module C* runs CVC on the SMT file and if the negation is satisfiable, interprets the counter-model in natural language.

shares the fundamental idea of employing LLMs to construct proofs and analyze relationships between textual spans, our task adds a layer of contextual reasoning by requiring the incorporation of background knowledge and maintaining interdependency between proof steps, which is not present in approaches where each proof step is treated as an independent, isolated process.

Theory solvers. Recent work by Hahn et al. (2022) demonstrated the potential of integrating symbolic solvers with large language models (LLMs), such as tool-augmented LLMs, to combine neural and symbolic reasoning. While such approaches are promising, they often struggle to translate natural language into symbolic representations and effectively capture background knowledge. Other recent approaches (Olausson et al., 2023; Pan et al., 2023) have used theory solvers to logically reason with natural language, which we build on with several key advancements. First, we introduce a framework that handles naturalistic, real-world data and tasks with ambiguous premises and conclusions. Then, we present a method to incorporate background knowledge into logical formulas. Finally, we show that our approach introduces interpretability by allowing human verification and modification throughout the intermediate reasoning steps.

3 Methodology

Although powerful, LLMs struggle to detect logical fallacies in language, as it requires proper logical analysis (Jin et al., 2022). On the other hand, SMT solvers can reason over logical formulas with theoretical guarantees but require the input to be in a structured, logical form. This approach combines

the strengths of both to classify logical fallacies.

Task formulation. The task input is an argument in natural language comprising one or more sentences, which is converted into formal logical form using a chain of LLMs. Following this, an SMT solver processes the logical form and returns whether it is valid. If invalid, the SMT solver provides a counterexample explaining why it is a logical fallacy, which is then interpreted with an LLM.

First-order logic. In FOL, propositions are represented using predicates that express properties or relations over objects in a domain. These predicates can be combined with constants, representing specific objects and variables that represent unspecified elements in the domain. An Interpretation assigns meaning to these symbols within a given context, while a Sort categorizes objects into different types, facilitating precise reasoning about their properties. Logical connectives of FOL, such as implication (\Rightarrow), universal quantifiers (\forall), existential quantifiers (\exists), and operators for conjunction/and (\wedge), disjunction/or (\vee), and negation/not (\neg), allow for the construction of intricate statements.

Module A: Natural language to first-order logic. Our approach for converting given natural language sentences into a logical form comprises multiple steps involving few-shot prompting of LLMs: (i) decomposing a sentence into multiple smaller parts that can be represented in first-order logic, (ii) identifying relationships between different sub-components to merge them and obtain a resultant logical formula, and (iii) identifying real-world relationships between these sub-components (background knowledge) and augmenting them to ob-

tain a FOL formula by incorporating background knowledge in the statement. We demonstrate with a Logical Fallacy (LF) and a Valid (V) example.

1. LF Example: A Logical Fallacy Input

I met a tall man who loved to eat cheese, now I believe all tall people like cheese.

2. V Example: A Valid Input

A boy is jumping on a skateboard in the middle of a red bridge. Thus the boy does a skateboarding trick.

Our pipeline begins with a semantic decomposition module which decomposes natural language arguments into respective claims and implications. Generally, a sentence can be split into some claims and implications based on those claims (see Prompt 2).

1. LF Example: Claim and Implication Parser

Claim: A tall man loved to eat cheese.
Implication: All tall people like cheese.

2. V Example: Claim and Implication Parser

Claim: A boy is jumping on a skateboard in the middle of a red bridge.
Implication: The boy does a skateboarding trick.

The claims and implications are split into further sub-components and used to build up the logical form of the sentence. The next step is to identify entities in the sentence. In our work, we treat noun phrases or surrogates for noun phrases as entities (see Prompt 3). Then, we find the relationship between the different entities using Zero-Shot classification via Natural Language Inference (NLI). These relationships (e.g., subset, equality, not related) are generally helpful in deciding appropriate quantifiers in the logical form. For example, if the entities are *man* and *people*, then it can be inferred that *man* is a subset of *people* and that the man would be bound by an existential quantifier in the sentence x (see Prompt 4).

1. LF Example: Entity Extractor

Referring expressions:

- man: x
- cheese: c
- people: y
- $x \subseteq y$

2. V Example: Entity Extractor

Referring expressions:

- boy: b
- skateboard: s
- bridge
- skateboardingTrick: y

The other set of sub-components are properties, which describe a trait of a referring expression or relationship between multiple referring expressions. These properties are predicates in first-order logic. We use a single module to extract the properties and

the relation between properties and entities. (see Prompt 5). We also find the relationships between various properties (see Prompt 6). For instance, in the LF Example, it can be inferred that *Like* and *Love* are contextually similar. Similarly, in our valid example, *jumping over skateboard* implies *doing a skateboard trick*. These relationships provide an additional context that is not directly present in the statement.

To identify these contextual relationships, we run NLI between each pair of properties, i.e., by setting one property as the hypothesis and the other as the premise as the input to the NLI model. If we find that any one property entails the other, we add the relationship $\text{property1} \Rightarrow \text{property2}$ to our context. Before running the NLI model between a pair of properties, we replace the variables in each property with the referring expressions that they represent. This adds additional context that helps the NLI model identify relations. For instance, in the V Example, the NLI model is unable to find the relation between *JumpsOn*(x, s) and *Does*(x, y), but it can identify the relationship between *JumpsOn*(boy, skateboard) and *Does*(boy, skateboardingTrick).

1. LF Example: Property Extractor + Background Knowledge Retriever

Properties: Tall, Love, Like

Property entity relations: Tall(x), Love(x, c)

Background knowledge:

1. $\forall x(\text{Like}(x, c) \Rightarrow \text{Love}(x, c))$
2. $\forall x(\text{Love}(x, c) \Rightarrow \text{Like}(x, c))$
3. $x \subseteq y$

2. V Example: Property Extractor + Background Knowledge Retriever

Properties: JumpsOn, inMiddleOf, Red, Does

Property entity relations: JumpsOn(b, s),

Red(bridge), inMiddleOf(b , bridge), Does(b, y)

Background knowledge:

1. $\forall x(\text{JumpsOn}(b, s) \Rightarrow \text{Does}(b, y))$

Finally, we combine all of this information using the relationships between properties and entities to obtain the FOL form of the sentence with the help of an LLM (see Prompt 7). For a logical fallacy, the negation of the formula is expected to be satisfiable. On the contrary, for a valid statement, the negation of the formula should be unsatisfiable.

1. LF Example: NL2FOL Output

First-order logic: $((\forall x(\text{Like}(x, c) \Rightarrow \text{Love}(x, c))) \wedge (\forall x(\text{Love}(x, c) \Rightarrow \text{Like}(x, c))) \wedge (\exists x(\text{Tall}(x) \wedge \text{Love}(x, c)))) \Rightarrow (\forall y(\text{Tall}(y) \Rightarrow \text{Like}(y, c))))$

2. V Example: NL2FOL Output

First-order logic: $(\forall x(\text{JumpsOn}(x, s) \Rightarrow \text{Does}(x, y)) \wedge \text{Red}(\text{bridge}) \wedge \text{inMiddleOf}(b, \text{bridge}) \wedge \text{JumpsOn}(b, s)) \Rightarrow \text{Does}(b, y)$

Module B: First-order logic to SMT. The next step involves automatically creating an SMT file for the negation of the first-order logical formula generated. While one can easily write an SMT file for a logical formula manually, generating one automatically for an arbitrary formula has not been done before. Thus, we develop a compiler that parses a given logical formula and converts it into an SMT file that can be given to CVC as input, as described in Algorithm 1 (See Appendix).

Module C: Interpreting SMT results. To verify the validity of the logical formulas, we utilize an SMT solver, CVC4 (Barrett et al., 2011). The solver determines whether the formula is valid or invalid, hence a logical fallacy. In the case of invalidity, the model provides a counterexample to the original logical formula, which shows that the given claim or implication is a logical fallacy.

Example (Module B Output):

I met a tall man who loved to eat cheese, now I believe all tall people like cheese.

↓
First-order logic: $((\forall x(\text{Like}(x, c) \Rightarrow \text{Love}(x, c))) \wedge (\forall x(\text{Love}(x, c) \Rightarrow \text{Like}(x, c))) \wedge (\exists x(\text{Tall}(x) \wedge \text{Love}(x, c)))) \Rightarrow (\forall y(\text{Tall}(y) \Rightarrow \text{Like}(y, c)))$

↓
SMT classification: Logical fallacy
Explanation: Counterexample

- ↓
- John is tall ($\text{Tall}(\text{John})$ is True). John likes cheese ($\text{Likes}(\text{John}, \text{Cheese})$ is True).
 - Jane is tall ($\text{Tall}(\text{Jane})$ is True). No constraint Jane likes cheese.

Therefore, there exists a tall person (John) who likes cheese, but it does not follow that all tall people like cheese, since Jane serves as a counterexample.

Figure 2: Example of logical fallacy detection using NL2FOL. The resulting classification is explained using a counterexample generated by the SMT solver.

The result of the SMT solver is hard to interpret, as it uses technical terminology generally only well understood by those who are familiar with CVC4 and SMT. To obtain an explanation in natural language, we prompt an LLM with the claim, implication, referring expressions, properties, FOL formula, and the counterexample generated by CVC4. The model then interprets the counterexample with natural language, as depicted in Figure 2.

4 Experiments

We evaluate our approach on both logical fallacies (positive class) and valid statements (negative

class). For logical fallacies, we use the LOGIC and LOGICCLIMATE (Jin et al., 2022) datasets, originally designed for training models to identify and classify different fallacies. These datasets contain examples of logical fallacies, each labeled with multiple categories from 13 different categories, including faulty generalization, circular claim, and ad hominem. The LOGIC dataset contains 2,449 examples of common logical fallacies collected mostly from quiz websites. The LOGICCLIMATE dataset comprises 1,079 examples of logical fallacies drawn from climate change news articles on the Climate Feedback platform. It is intended to test the model’s ability to generalize out-of-domain.

To test our approach with valid statements, we use the Stanford Natural Language Inference (SNLI) corpus (Bowman et al., 2015), which supports the development of natural language inference systems. This dataset features over 570,000 human-annotated sentence pairs, where each pair consists of a premise and a hypothesis labeled as entailment, contradiction, or neutral. We focus on the entailment class in this study, extracting over 170,000 sentence pairs where the premise entails the hypothesis. We construct valid sentences by combining the premise and hypothesis into a single sentence.

The task is set up as a simple binary classification task, where the input consists of sentences drawn from the LOGIC or LOGICCLIMATE datasets labeled as logical fallacies or from the SNLI dataset labeled as valid sentences. Here, we treat logical fallacies as the positive class. To ensure a balanced evaluation, we select an equal number of fallacies and valid statements, allowing for a fair comparison across both classes. Finally, our model is evaluated on standard binary classification metrics such as precision, recall, f1 score, and accuracy.

Models. We compare our method to pretrained language models, including Llama2-7B (Touvron et al., 2023), GPT4o-mini (OpenAI, 2024), GPT4o (OpenAI et al., 2024a) and OpenAI o1-preview (OpenAI et al., 2024b) with few-shot in-context examples (see Prompt 1). We also run NL2FOL with each of the above models used for the LLM prompting stages. Llama2-7B was chosen for our experiments as it had the best performance during testing over an initial subset of the data, outperforming Llama3.1-8B (Grattafiori et al., 2024), Llama3.2-11B (AI, 2024a), and Ministral-8B (AI, 2024b). We evaluate BART (140M parameters)

Model	Method	LOGIC				LOGICCLIMATE			
		Acc.	P.	R.	F1	Acc.	P.	R.	F1
Llama-7B	End-to-end	0.41	0.45	0.82	0.58	0.31	0.38	0.62	0.47
	NL2FOL (Ours)	0.63	0.58	0.92	0.71	0.66	0.60	0.94	0.73
GPT-4o-mini	End-to-end	0.91	0.94	0.88	0.91	0.64	0.67	0.55	0.60
	NL2FOL (Ours)	0.70	0.64	0.91	0.75	0.73	0.66	0.93	0.77
GPT-4o	End-to-end	0.96	0.96	0.96	0.96	0.70	0.95	0.42	0.58
	NL2FOL (Ours)	0.78	0.76	0.82	0.78	0.80	0.80	0.80	0.80
OpenAI o1-preview	End-to-end	0.93	0.89	0.98	0.93	0.73	0.84	0.56	0.67
	NL2FOL (Ours)	-	-	-	-	-	-	-	-

Table 2: Comparison of few-shot model performance metrics (abbreviations: Acc. = accuracy, P. = precision, R. = recall, F1 = F1 score) on the LOGIC+SNLI and LOGICCLIMATE+SNLI datasets using End-to-end vs. NL2FOL (Ours). Results on NL2FOL with o1-preview are omitted as o1-preview failed to complete the pipeline in most cases, likely due to its poor instruction following capabilities.

(Lewis et al., 2020) finetuned on MNLI (Williams et al., 2018) to analyze the relationships between properties and referring expressions. We ran the experiments on a V100 GPU, with one run costing around 2 GPU hours.

Prompt tuning. For prompt tuning, 20 samples from the LOGIC dataset were selected and manually annotated with intermediate and final results. They were then split into 10 train and 10 validation examples. For each prompt, we start with a simple description of the task. 4-6 examples were randomly selected from the train set as in-context examples, with the relevant intermediate outputs depending on the stage. Results were tested on the validation examples, and the prompt was updated to address common mistakes. To ensure fairness, a fixed number of 5 improvement iterations was used for each prompt, and the one showing best performance over the validation examples was chosen.

5 Results and Discussion

As shown in Table 2, our method achieves an F1 score of 78% when used with GPT-4o on the LOGIC dataset. When run end-to-end, the Llama-7B model reached an F1 score of only 58%, but when used with the NL2FOL pipeline, reached a score of 71%. Although end-to-end classification has shown better performance in other models, comparisons can be skewed because they may have been exposed to the LOGIC dataset and its labels during training because this dataset was compiled from publicly accessible web sources. On average, NL2FOL demonstrated high recall, whereas end-to-end classification demonstrated high precision.

Our challenge set LOGICCLIMATE+SNLI contains real-world logical fallacies from climate change

news. Since this dataset was used to test generalization, the in-context examples we provide to all models are from the LOGIC dataset. NL2FOL yields results that are highly similar to the results from LOGIC, whereas end-to-end classification saw a drop in performance. This demonstrates that our system is also robust and adapts well to real-world texts, including texts with significant domain-specific context. This makes it effective in detecting and mitigating misinformation. Specifically, on this dataset, we find that NL2FOL outperforms direct translation with all LLMs that we tested.

5.1 Quantitative Analysis

Error analysis and interpretability. The proposed method is interpretable due to the use of natural language inputs and outputs at each step of the pipeline. This structure allows for precise identification of the specific module responsible for a failure by examining intermediate results. To evaluate this aspect, we performed an in-depth error analysis by annotating the module responsible for failure in 100 incorrect predictions made by the model. The results are summarized in Table 4.

Our analysis reveals that the majority of errors occur in the ‘Background Knowledge Retriever’, involving missed or incorrectly added contextual information in the logical form. Other errors typically pertain to incorrect identification of claims, implications, or properties. In contrast, inaccuracies in the generation of logical forms are relatively infrequent, suggesting that the model performs well in constructing accurate logical representations when provided with reliable information about the constituent entities and properties within a sentence. This finding underscores the importance of improving the background knowledge retriever module to

Type	Sentence	Logical Form	Prediction
1 LF	X has been around for years now. Y is new. Therefore, Y is better than X.	$(\text{IsNew}(a) \wedge \sim \text{IsNew}(b)) \Rightarrow (\text{IsBetterThan}(a,b))$	LF: Correct prediction
2 LF	Everyone is doing the Low-Carb Diet.	$(\exists b (\exists a (\text{IsDoing}(b,a)))) \Rightarrow (\exists c (\exists a (\text{IsDoing}(c,a))))$	V: Incorrect prediction - Wrong translation given when no claim given
3 V	Two dogs are fighting in a field. Consequently, the two dogs are outside.	$(\exists b (\exists a (\text{IsFighting}(a, b) \wedge \text{IsInField}(b) \wedge \text{IsInField}(b)))) \Rightarrow (\exists a (\text{IsOutside}(a)))$	LF: Incorrect prediction - Missing semantic ground truth claim: $\forall a (\text{IsInField}(a) \Rightarrow \text{IsOutside}(a))$
4 V	A baseball player gets ready to catch a fly ball near the outfield fence. Therefore, a person is playing baseball outdoors.	$(\exists a (\text{IsGettingReady}(a) \wedge (\text{IsABaseballPlayer}(a) \wedge \text{IsCatchingFlyBall}(a) \wedge \text{IsNearOutfieldFence}(a)))) \wedge (\forall e (\text{IsABaseballPlayer}(e) \Rightarrow \text{IsPlayingBaseball}(e))) \wedge (\forall f (\text{IsPlayingBaseball}(f) \Rightarrow \text{IsABaseballPlayer}(f))) \wedge (\forall g (\text{IsNearOutfieldFence}(g) \Rightarrow \text{IsOutdoors}(g)))) \Rightarrow (\exists c (\exists a (\text{IsPlayingBaseball}(a) \wedge \text{IsOutdoors}(c))))$	V: Correct Prediction - The method identifies additional context by establishing relationships such as <i>IsBaseballPlayer</i> implying <i>IsPlayingBaseball</i> , and <i>IsNearOutfieldFence</i> implying <i>IsOutdoors</i> .
5 V	A woman sits alone on a park bench in the sun. Hence, a woman is in a park.	$(\text{IsSittingOn}(a, b) \wedge \text{isParkBench}(b) \wedge \text{IsInSun}(a)) \Rightarrow (\text{IsInPark}(a))$	LF: Incorrect prediction - Missing semantic ground truth claim: $\forall a \forall b (\text{IsSittingOn}(a, b) \wedge \text{isParkBench}(b) \Rightarrow \text{IsInPark}(a))$
6 V	A woman is standing at a podium. Thus, a person is at a podium.	$(\exists a \exists b (\text{IsStandingAt}(b, a)) \wedge \forall f \forall e \forall d (\text{IsStandingAt}(d,e) \Rightarrow \text{IsAt}(f,e)) \Rightarrow \exists c \exists a (\text{IsAt}(c, a))$	V: Correct prediction - The method identifies additional context by establishing the relationship <i>IsStandingAt</i> implying <i>IsAt</i> .

Table 3: Some example outputs of our model (abbreviations: LF = Logical Fallacy, V = Valid statement)

Sub-Module with Error	Error Proportion
Claim and Implication Parser	0.19
Incorrect Label	0.01
Property Extractor	0.13
Background Knowledge Retriever	0.54
FOL Formulation Engine	0.13

Table 4: Categorization of model errors by type on NL2FOL (GPT-4o), based on a review by domain experts in the logic of 100 randomly sampled examples

improve overall model performance.

Impact of adding background knowledge to NL2FOL. Based on the error analysis, missing or incorrect background knowledge was a significant contributor to incorrect predictions of our method. To quantitatively assess the impact of grounding on model performance, we evaluated several approaches for NLI in the Background Relation Extractor. These included: (a) a pipeline without any background knowledge as a baseline, (b) a model without context where the LLM (GPT4o) only processes the input properties, (c) an LLM that incorporates both the input sentence and properties and (d) a smaller model specifically fine-tuned for NLI (BART-MNLI). Results are presented in Table 5.

We see that precision and recall both improve sig-

nificantly with better grounding techniques. The LLM model with sentence context achieves the highest overall performance. This is likely due to the sentence context providing information about clauses that are omitted due to the choice of representation in FOL. This indicates that integrating robust grounding mechanisms is critical to enhancing the accuracy and reliability of the method.

Method	LOGIC+SNLI				LOGICCLIMATE+SNLI			
	Acc.	P.	R.	F1	Acc.	P.	R.	F1
(a) No Grounding	0.54	0.52	0.88	0.66	0.57	0.54	0.94	0.69
(b) LLM	0.76	0.78	0.74	0.75	0.79	0.80	0.78	0.79
(c) LLM w/ context	0.78	0.76	0.82	0.78	0.80	0.80	0.80	0.80
(d) BART-MNLI	0.71	0.71	0.70	0.70	0.77	0.81	0.71	0.77

Table 5: Comparison of different grounding methods on NL2FOL (GPT4o-mini) across the LOGIC+SNLI and LogicClimate+SNLI datasets

Impact of using an SMT solver. To assess the impact of using an SMT solver in our pipeline, we compared its performance against an LLM as a baseline for classifying the logical forms as valid or fallacies. The results, summarized in Table 6, demonstrate a significant improvement in performance metrics with the integration of the SMT solver. Results reveal the SMT-based approach significantly outperforms the LLM-based approach in all metrics across both the LOGIC and LOGIC-

CLIMATE datasets. This underscores the advantage of formal reasoning systems like SMT solvers for tasks requiring precise logical inference and structured reasoning compared to LLMs, which may lack systematic consistency in such contexts.

5.2 Qualitative Analysis

5.2.1 Success Modes of NL2FOL

S1: Captures implicit information not mentioned in premises. Previous works that directly translate natural language to logical forms suffer from an inability to capture implicit information not mentioned in the premises (Olausson et al., 2023). Our ‘Background Knowledge Retriever’ step allows us to capture this information in the final logical form. An illustration of this can be found in Example 4 of Table 3.

S2: Captures explicit information that is missed in the representation. Our pipeline is also able to capture information that is explicitly mentioned in the premises but missed due to the choice of representation in logical form. In Example 6, in Table 3, the fact that the woman is both standing and is at the podium is lost due to the choice representation *IsStandingAt*. However, the fact that the woman is at the podium is recovered in the final logical form due to the identified background knowledge *IsStandingAt* implies *IsAt*.

S3: Comparison to direct translation. To evaluate the efficacy of the multi-step LLM pipeline, we compared it against a direct translation approach, where natural language inputs were converted into logical forms with a single LLM call using a few-shot prompt. However, this task proved to be excessively complex for LLMs. Llama failed to generate any output, citing an inability to comprehend the prompt. Larger LLMs exhibited significant limitations, with over 95% of their outputs containing syntax errors. These findings highlight the inadequacy of direct translation for complex logical reasoning tasks and underscore the necessity of a structured, multi-step approach to ensure the accuracy and syntactic correctness of the logical form.

5.2.2 Failure Modes of NL2FOL

F1: Misses some background knowledge. As can be observed in Table 4, incorrect identification of background knowledge is the most common cause for incorrect classifications. This is because any gaps in background knowledge can cause a valid statement to be identified as a logical fallacy, and

Classifier	LOGIC				LOGICCLIMATE			
	Acc.	P.	R.	F1	Acc.	P.	R.	F1
SMT	0.78	0.76	0.82	0.78	0.80	0.80	0.80	0.80
GPT-4o	0.69	0.71	0.62	0.66	0.73	0.72	0.74	0.73

Table 6: Comparison of classification methods used with NL2FOL (GPT4o) on LOGIC and LOGICCLIMATE

an incorrectly added clause can cause a fallacy to be identified as valid. One such case is present in example 3 of the Table 3. In this case, the model is not able to identify the extra context statement because the NLI model does not identify a required ground-truth relation. If this context were to be added to the claim of the logical formula, then the statement would have been predicted to be valid.

F2: Limitations of NLI. Our current approach is limited to discerning relationships between two properties at a time rather than handling multiple relationships concurrently. For reference, consider Example 5 in Table 3. Here, the semantic claim involves the conjunction of two properties entailing the third, while the ‘Background Knowledge Retriever’ only checks whether one property entails the other. Finding such complex extra context requires more advanced techniques or additional human intervention. Including them could further improve the precision of the model overall.

F3: Imprecision of LLMs. Among the logical fallacies that our model incorrectly predicted to be a valid statement, most of these predictions failed due to the imprecision of the LLM, leading to false translations and incorrect results. Example 2 demonstrates a case where the input does not have any claim but instead jumps straight to an implication. However, the model is not able to identify that the example has no claim. As a result, we obtain an incorrect translation with our technique.

6 Conclusion

We present an effective and automatic solution to detect fallacies and tackle misinformation. We developed a strategy to distinguish logical fallacies from valid statements, involving a chaining approach to convert a sentence to first-order logic using LLMs, followed by using SMT solvers to identify whether the first-order logical statement is valid or not. If not, we interpret the counter-model generated by the SMT solver in natural language. Our proposed technique shows promising results in identifying logical fallacies and valid statements, as well as good generalizability across domains.

Ethics Statement

While the intended outcome of this research is to help fight misinformation and promote rational discourse, there are several ethical challenges that we must consider. First, dependence on AI to identify logical fallacies could influence how individuals engage in debates and discussions. There is a risk that people may over-rely on AI judgments, potentially stifling complex statements or dissenting opinions that are essential for a healthy democratic process. Moreover, the use of AI in moderating discussions, especially in identifying logical fallacies, raises ethical questions about the automation of content moderation. While it can enhance the quality of public discourse by filtering out fallacious statements, it also risks automating censorship and impacting the dynamics of online communities. In the wrong hands, logical fallacy detection tools could be exploited to silence speech or suppress viewpoints under the pretext of promoting rational discourse. This potentially allows governments or organizations to stifle opposition or critique.

To address these issues, we advocate for the development of ethical guidelines for AI use that emphasize transparency, accountability, and active user engagement. These measures are crucial in encouraging public literacy in AI and logical fallacies, ultimately empowering individuals to critically assess both AI output and arguments they may encounter.

Limitations

Scope of logical reasoning tasks. Correct identification of background knowledge is crucial for our method. While we have shown its potential in detecting logical fallacies for short and structured premises, it is important to note that this approach may miss complex relational constructs (for example, $(a \wedge b) \Rightarrow (c \vee d)$), in which richer logical patterns may often be required in real-world reasoning tasks such as those present in multi-paragraph contexts or Question-Answering (QA) datasets.

Generalizability to other tasks and domains. We have demonstrated promising results of our approach to logical fallacy detection, but whether the findings generalize to other logical tasks and domains remains unexplored. The performance of our approach in other languages is untested and may introduce unforeseen challenges.

Going beyond first-order logic. It is unknown

whether our approach would be sufficiently expressive for reasoning tasks requiring higher-order or non-classical logic, as we limit our exploration to first-order logic. Conceptually, extending our method to the aforementioned domains is feasible but would require modification to the SMT integration and LLM-driven logic translation processes. Thus, further testing may include translating to logic beyond FOL, such as temporal and higher-order logic.

Computational cost. Using LLMs and SMT solvers can incur high computational costs, such as high-performance GPUs for LLM inference, CPUs optimized for SMT solvers, and high API usage, particularly for models like GPT-o1 and Llama-7B.

References

- Meta AI. 2024a. [Llama 3.2-11b model card](#). Accessed: 2025-02-15.
- Mistral AI. 2024b. [Ministral-8b-instruct-2410 model card](#). Accessed: 2025-02-15.
- Tariq Alhindi, Tuhin Chakrabarty, Elena Musi, and Smaranda Muresan. 2022. [Multitask instruction-based prompting for fallacy recognition](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8172–8187, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Rami Aly, Marek Strong, and Andreas Vlachos. 2023. [QA-natver: Question answering for natural logic-based fact verification](#). In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Gabor Angeli and Christopher D Manning. 2014. Natu-ralli: Natural logic inference for common sense reasoning. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 534–545.
- Bing Bai, Jian Liang, Guanhua Zhang, Hao Li, Kun Bai, and Fei Wang. 2020. [Why attentions may not be interpretable?](#) *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*.
- Clark Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanović, Tim King, Andrew Reynolds, and Cesare Tinelli. 2011. Cvc4. In *Computer Aided Verification*, pages 171–177, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Clark Barrett, Aaron Stump, and Cesare Tinelli. 2009. Satisfiability modulo theories. *Communications of the ACM*, 52(9):69–77.
- Iz Beltagy, Stephen Roller, Pengxiang Cheng, Katrin Erk, and Raymond J. Mooney. 2016. Representing meaning with a combination of logical and distributional models. *Computational Linguistics*, 42(4):763–808.

666	Samuel R. Bowman, Gabor Angeli, Christopher Potts,	Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, An-	723
667	and Christopher D. Manning. 2015. A large annotated	thony Hartshorn, Aobo Yang, Archi Mitra, Archie Sra-	724
668	corpus for learning natural language inference . In <i>Pro-</i>	vankumar, Artem Korenev, Arthur Hinsvark, Arun Rao,	725
669	<i>ceedings of the 2015 Conference on Empirical Methods</i>	Aston Zhang, Aurelien Rodriguez, Austen Gregerson,	726
670	<i>in Natural Language Processing</i> , pages 632–642, Lis-	Ava Spataru, Baptiste Roziere, Bethany Biron, Binh	727
671	bon, Portugal. Association for Computational Linguis-	Tang, Bobbie Chern, Charlotte Caucheteux, Chaya	728
672	tics.	Nayak, Chloe Bi, Chris Marra, Chris McConnell,	729
673	Tom B Brown, Benjamin Mann, Nick Ryder, Melanie	Christian Keller, Christophe Touret, Chunyang Wu,	730
674	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Corinne Wong, Cristian Canton Ferrer, Cyrus Niko-	731
675	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	laidis, Damien Allonsius, Daniel Song, Danielle Pintz,	732
676	Askeel, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen	Danny Livshits, Danny Wyatt, David Esiobu, Dhruv	733
677	Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,	Choudhary, Dhruv Mahajan, Diego Garcia-Olano,	734
678	Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris	Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab	735
679	Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott	AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael	736
680	Gray, Benjamin Chess, Jack Clark, Christopher Berner,	Smith, Filip Radenovic, Francisco Guzmán, Frank	737
681	Sam McCandlish, Alec Radford, Ilya Sutskever, and	Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis	738
682	Dario Amodei. 2020. Language models are few-shot	Anderson, Govind Thattai, Graeme Nail, Gregoire Mi-	739
683	learners. <i>Advances in Neural Information Processing</i>	alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,	740
684	<i>Systems</i> , 33:1877–1901.	Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	741
685	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,	Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra,	742
686	Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee,	Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan	743
687	Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori,	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet	744
688	Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang.	Shah, Jelmer van der Linde, Jennifer Billock, Jenny	745
689	2023. Sparks of artificial general intelligence: Early	Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu	746
690	experiments with gpt-4 .	Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bit-	747
691	Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji,	ton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua	748
692	and Quanquan Gu. 2024. Self-play fine-tuning converts	Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden	749
693	weak language models to strong language models .	Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plaw-	750
694	Matthias Cosler, Christopher Hahn, Daniel Mendoza,	iak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid	751
695	Frederik Schmitt, and Caroline Trippel. 2023. nl2spec:	El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu,	752
696	Interactively translating unstructured natural language	Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Young,	753
697	to temporal logics with large language models . In <i>Com-</i>	Laurens van der Maaten, Lawrence Chen, Liang Tan,	754
698	<i>puter Aided Verification. CAV 2023. Lecture Notes in</i>	Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,	755
699	<i>Computer Science</i> , volume 13965, Cham. Springer.	Lukas Blecher, Lukas Landzaat, Luke de Oliveira,	756
700	Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zheng-	Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh,	757
701	nan Xie, Hannah Smith, Leighanna Pipatanangkura, and	Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli,	758
702	Peter Clark. 2021. Explaining answers with entailment	Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie	759
703	trees . In <i>Proceedings of the 2021 Conference on Em-</i>	Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh,	760
704	<i>pirical Methods in Natural Language Processing</i> , pages	Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay	761
705	7358–7370, Online and Punta Cana, Dominican Repub-	Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning	762
706	lic. Association for Computational Linguistics.	Zhang, Olivier Duchenne, Onur Çelebi, Patrick Al-	763
707	Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: an	rassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter	764
708	efficient smt solver. In <i>Proceedings of the Theory and</i>	Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krish-	765
709	<i>Practice of Software, 14th International Conference on</i>	nan, Punit Singh Koura, Puxin Xu, Qing He, Qingx-	766
710	<i>Tools and Algorithms for the Construction and Analy-</i>	iao Dong, Ragavan Srinivasan, Raj Ganapathy, Ra-	767
711	<i>sis of Systems, TACAS’08/ETAPS’08</i> , page 337–340,	mon Calderer, Ricardo Silveira Cabral, Robert Stojnic,	768
712	Berlin, Heidelberg. Springer-Verlag.	Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Ro-	769
713	Francesco Fuggitti and Tathagata Chakraborti. 2023.	hit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan	770
714	Nl2l - a python package for converting natural lan-	Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang,	771
715	guage (nl) instructions to linear temporal logic (ltl) for-	Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh,	772
716	mulas . In <i>AAAI Conference on Artificial Intelligence</i> .	Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shao-	773
717	Gaël Gendron, Qiming Bao, Michael Witbrock, and	liang Nie, Sharan Narang, Sharath Raparthy, Sheng	774
718	Gillian Dobbie. 2024. Large language models are not	Shen, Shengye Wan, Shruti Bhosale, Shun Zhang,	775
719	strong abstract reasoners .	Simon Vandenhende, Soumya Batra, Spencer Whit-	776
720	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	man, Sten Sootla, Stephane Collot, Suchin Gururan-	777
721	Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,	gan, Sydney Borodinsky, Tamar Herman, Tara Fowler,	778
722	Aiesha Letman, Akhil Mathur, Alan Schelten, Alex	Tarek Sheasha, Thomas Georgiou, Thomas Scialom, To-	779
		bias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	780
		Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ra-	781
		manathan, Viktor Kerkez, Vincent Gouget, Virginie	782
		Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Wei-	783
		wei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers,	784
		Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xi-	785
		aoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia,	786

787	Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yas-	Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan	851
788	mine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue	Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,	852
789	Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng	Piotr Dollar, Polina Zvyagina, Prashant Ratanchan-	853
790	Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh,	dani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	854
791	Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	855
792	Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Ray-	856
793	Goldstand, Ajay Menon, Ajay Sharma, Alex Boesen-	mond Li, Rebekkah Hogan, Robin Battey, Rocky Wang,	857
794	berg, Alexei Baevski, Allie Feinstein, Amanda Kallet,	Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby,	858
795	Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu,	Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara	859
796	Andres Alvarado, Andrew Caples, Andrew Gu, Andrew	Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan,	860
797	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto,	861
798	dani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita	Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-	862
799	Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,	863
800	Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan,	Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang	864
801	Beau James, Ben Maurer, Benjamin Leonhardi, Bernie	Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe,	865
802	Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape,	Soumith Chintala, Stephanie Max, Stephen Chen, Steve	866
803	Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram	Kehoe, Steve Satterfield, Sudarshan Govindaprasad,	867
804	Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido,	Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk,	868
805	Britt Montalvo, Carl Parker, Carly Burton, Catalina	Suraj Subramanian, Sy Choudhury, Sydney Goldman,	869
806	Mejia, Ce Liu, Changan Wang, Changkyu Kim, Chao	Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	870
807	Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	871
808	Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon	Matthews, Timothy Chou, Tzook Shaked, Varun Von-	872
809	Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David	timitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,	873
810	Adkins, David Xu, Davide Testuggine, Delia David,	Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad	874
811	Devi Parikh, Diana Liskovich, Didem Foss, Dingkan	Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov,	875
812	Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa	Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,	876
813	Jamil, Elaine Montgomery, Eleonora Presani, Emily	Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan	877
814	Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman,	Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun	878
815	Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei	Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying	879
816	Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat	Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao,	880
817	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait,	881
818	Seide, Gabriella Medina Florez, Gabriella Schwarz,	Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu	882
819	Gada Badeer, Georgia Swee, Gil Halpern, Grant Her-	Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3	883
820	man, Grigory Sizov, Guangyi, Zhang, Guna Lakshmi-	herd of models .	884
821	narayanan, Hakan Inan, Hamid Shojanazeri, Han Zou,		
822	Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison	Ivan Habernal, Henning Wachsmuth, Iryna Gurevych,	885
823	Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-	and Benno Kiesel. 2018a. "dummy, grandpa, do you	886
824	man, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog,	know anything?": Identifying and characterizing ad	887
825	Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai	hominem fallacies in the wild. In <i>Proceedings of the</i>	888
826	Gat, Jake Weissman, James Geboski, James Kohli, Jan-	<i>12th International AAAI Conference on Web and Social</i>	889
827	ice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Mar-	<i>Media (ICWSM)</i> , pages 206–215.	890
828	cus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy		
829	Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin,	Ivan Habernal, Henning Wachsmuth, Iryna Gurevych,	891
830	Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shep-	and Benno Stein. 2018b. Before name-calling: Dynam-	892
831	ard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg,	ics and triggers of ad hominem fallacies in web argu-	893
832	Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kar-	mentation . In <i>Proceedings of the 2018 Conference of</i>	894
833	tikay Khandelwal, Katayoun Zand, Kathy Matosich,	<i>the North American Chapter of the Association for Com-</i>	895
834	Kaushik Veeraraghavan, Kelly Michelena, Keqian Li,	<i>putational Linguistics: Human Language Technologies,</i>	896
835	Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	<i>Volume 1 (Long Papers)</i> , pages 386–396, New Orleans,	897
836	Huang, Lailin Chen, Lakshya Garg, Lavender A, Le-	Louisiana. Association for Computational Linguistics.	898
837	andro Silva, Lee Bell, Lei Zhang, Liangpeng Guo,		
838	Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Ma-	Christopher Hahn, Frederik Schmitt, Julia J. Tillman,	899
839	dian Khabsa, Manav Avalani, Manish Bhatt, Martyn	Niklas Metzger, Julian Siber, and Bernd Finkbeiner.	900
840	as Mankus, Matan Hasson, Matthew Lennie, Matthias	2022. Formal specifications from natural language .	901
841	Reso, Maxim Groshev, Maxim Naumov, Maya Lathi,		
842	Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal	Patrick M. Haluptzok, Matthew Bowers, and Adam Tau-	902
843	Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov,	man Kalai. 2022. Language modelxrls can teach them-	903
844	Mikayel Samvelyan, Mike Clark, Mike Macey, Mike	selves to program better . <i>ArXiv</i> , abs/2207.14502.	904
845	Wang, Miquel Jubert Hermoso, Mo Metanat, Moham-		
846	mad Rastegari, Munish Bansal, Nandhini Santhanam,	Zhijing Jin, Abhinav Lalwani, Tejas Vaidhya, Xiaoyu	905
847	Natascha Parks, Natasha White, Navyata Bawa, Nayan	Shen, Yiwen Ding, Zhiheng Lyu, Mrinmaya Sachan,	906
848	Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta,	Rada Mihalcea, and Bernhard Schoelkopf. 2022. Logi-	907
849	Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng,	cal fallacy detection . In <i>Findings of the Association for</i>	908
850	Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem	<i>Computational Linguistics: EMNLP 2022</i> , pages 7180–	909

7198, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	969
Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880. Online. Association for Computational Linguistics.	970
Jason Xinyu Liu, Ziyi Yang, Benjamin Schornstein, Sam Liang, Ifrah Idrees, Stefanie Tellex, and Ankit Shah. 2022. Lang2LTL: Translating natural language commands to temporal specification with large language models. In <i>CoRL Workshop on Language and Robot Learning</i> .	971
B. MacCartney and C. D. Manning. 2014. Natural logic and natural language inference . In H. Bunt, J. Bos, and S. Pulman, editors, <i>Computing Meaning: Volume 4</i> , pages 129–147. Springer Netherlands, Dordrecht.	972
Callistus Ireneous Nakpih and Simone Santini. 2020. Automated discovery of logical fallacies in legal argumentation . <i>International Journal of Artificial Intelligence & Applications</i> .	973
Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum, and Roger Levy. 2023. LINC: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5153–5176, Singapore. Association for Computational Linguistics.	974
OpenAI. 2024. Gpt-4o mini: Advancing cost-efficient intelligence . Accessed: 2025-02-15.	975
OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Lerner, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang,	976
Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichen, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Hariman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondrasiuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljube, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement,	977

1033	Owen Campbell-Moore, Patrick Chao, Paul McMillan,	Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang,	1096
1034	Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Pe-	Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo,	1097
1035	ter Deng, Peter Dolan, Peter Hoeschele, Peter Welin-	Joel Parish, Johannes Heidecke, John Hallman, John	1098
1036	der, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla	Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan	1099
1037	Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim,	Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao,	1100
1038	Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo	Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi,	1101
1039	Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud	Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg,	1102
1040	Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob	Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama	1103
1041	Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-	Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Ley-	1104
1042	dani, Romain Huet, Rory Carmichael, Rowan Zellers,	ton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay	1105
1043	Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu,	McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kon-	1106
1044	Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer,	draciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd,	1107
1045	Samuel Miserendino, Sandhini Agarwal, Sara Culver,	Maja Trebacz, Manas Joglekar, Mark Chen, Marko	1108
1046	Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger,	Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max	1109
1047	Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sher-	Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y.	1110
1048	win Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia	Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mi-	1111
1049	Phene, Spencer Papay, Srinivas Narayanan, Steve Cof-	anna Chen, Michael Lampe, Michael Malek, Michele	1112
1050	fey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda,	Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov,	1113
1051	Tal Stramer, Tao Xu, Tarun Gogineni, Taya Chris-	Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavar-	1114
1052	tianson, Ted Sanders, Tejal Patwardhan, Thomas Cun-	ian, Mostafa Rohaninejad, Nat McAleese, Neil Chowd-	1115
1053	ninghman, Thomas Degry, Thomas Dimson, Thomas	hury, Neil Chowdhury, Nick Ryder, Nikolas Tezak,	1116
1054	Raoux, Thomas Shadwell, Tianhao Zheng, Todd Un-	Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk,	1117
1055	derwood, Todor Markov, Toki Sherbakov, Tom Rubin,	Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Iz-	1118
1056	Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Pe-	mailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Ran-	1119
1057	tersson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit	dall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara,	1120
1058	Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko,	Reimar Leike, Renny Hwang, Rhythm Garg, Robin	1121
1059	Wayne Chang, Weiye Zheng, Wenda Zhou, Wesam Man-	Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan	1122
1060	assra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei	Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam	1123
1061	Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen	Toyer, Samuel Miserendino, Sandhini Agarwal, Santi-	1124
1062	He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury	ago Hernandez, Sasha Baker, Scott McKinney, Scottie	1125
1063	Malkov. 2024a. Gpt-4o system card .	Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar,	1126
1064	OpenAI, Aaron Jaech, Adam Kalai, Adam Lerer, Adam	Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu,	1127
1065	Richardson, Ahmed El-Kishky, Aiden Low, Alec Hel-	Spencer Papay, Steph Lin, Suchir Balaji, Suvansh San-	1128
1066	lyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex	jeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang,	1129
1067	Iftimie, Alex Karpenko, Alex Tachard Passos, Alexan-	Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault	1130
1068	der Neitz, Alexander Prokofiev, Alexander Wei, Alli-	Sottiaux, Thomas Degry, Thomas Dimson, Tianhao	1131
1069	son Tam, Ally Bennett, Ananya Kumar, Andre Saraiva,	Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor	1132
1070	Andrea Vallone, Andrew Duberstein, Andrew Kon-	Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vi-	1133
1071	drich, Andrey Mishchenko, Andy Applebaum, An-	neet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad	1134
1072	gela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghor-	Fomenko, Weiye Zheng, Wenda Zhou, Wes McCabe,	1135
1073	bani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak,	Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining	1136
1074	Bob McGrew, Borys Minaiev, Botao Hao, Bowen	Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang,	1137
1075	Baker, Brandon Houghton, Brandon McKinzie, Brydon	Yunyun Wang, Zheng Shao, and Zhuohan Li. 2024b.	1138
1076	Eastman, Camillo Lugaresi, Cary Bassin, Cary Hud-	Openai o1 system card .	1139
1077	son, Chak Ming Li, Charles de Bourcy, Chelsea Voss,	Liangming Pan, Alon Albalak, Xinyi Wang, and	1140
1078	Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger,	William Wang. 2023. Logic-LM: Empowering large lan-	1141
1079	Christopher Hesse, Claudia Fischer, Clive Chan, Dan	guage models with symbolic solvers for faithful logical	1142
1080	Roberts, Daniel Kappler, Daniel Levy, Daniel Sel-	reasoning . In <i>Findings of the Association for Computa-</i>	1143
1081	sam, David Dohan, David Farhi, David Mely, David	<i>tional Linguistics: EMNLP 2023</i> , pages 3806–3824,	1144
1082	Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica,	Singapore. Association for Computational Linguistics.	1145
1083	Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth	William C Purdy. 1991. A logic for natural language.	1146
1084	Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik	<i>Notre Dame Journal of Formal Logic</i> , 32(3):409–425.	1147
1085	Ritter, Evan Mays, Fan Wang, Felipe Petroski Such,	Hrituraj Singh, Milan Aggarwal, and Balaji Krish-	1148
1086	Filippo Raso, Florencia Leoni, Foivos Tsimplouras,	namurthy. 2020. Exploring neural models for pars-	1149
1087	Francis Song, Fred von Lohmann, Freddie Sulit, Ge-	ing natural language into first-order logic . <i>ArXiv</i> ,	1150
1088	off Salmon, Giambattista Parascandolo, Gildas Chabot,	abs/2002.06544 .	1151
1089	Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi	Zhivar Sourati, Vishnu Priya Prasanna Venkatesh, Dar-	1152
1090	Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hes-	shan Deshpande, Himanshu Rawlani, Filip Ilievski,	1153
1091	sam Bagherinezhad, Hongyu Ren, Hunter Lightman,	Hông-Ân Sandlin, and Alain Mermoud. 2022. Robust	1154
1092	Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian	and explainable identification of logical fallacies in nat-	1155
1093	Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya		
1094	Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pa-		
1095	chocki, James Lennon, Jason Wei, Jean Harb, Jerry		

ural language arguments. *Knowledge Based Systems*, 266:110418.

Christian Stab and Iryna Gurevych. 2017. [Recognizing insufficiently supported arguments in argumentative essays](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 980–990, Valencia, Spain. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).

Jason Wei, Xuezhi Wang, Dale Schuurmans, et al. 2022. [Chain of thought prompting elicits reasoning in large language models](#). *Advances in Neural Information Processing Systems*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, and Faramarz Fekri. 2024. [Harnessing the power of large language models for natural language to first-order logic translation](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6942–6959, Bangkok, Thailand. Association for Computational Linguistics.

Appendix

A Algorithms

Algorithm 1: Compiling Logical Formula to SMT

Input: Logical formula \mathcal{L} in natural language or First-Order Logic (FOL)

Output: SMT file \mathcal{S} formatted for formal solvers

```
1 Step 1: Tokenize Formula
2  $\mathcal{T} \leftarrow \text{Tokenize}(\mathcal{L})$  // Split  $\mathcal{L}$  into tokens based
   on operators, parentheses, and commas
3 Step 2: Process Tokens
4  $\mathcal{P} \leftarrow \emptyset$  // Initialize processed tokens set
5 foreach token  $t \in \mathcal{T}$  do
6   if  $t$  is a predicate then
7     Identify arguments of  $t$ 
8     Recursively ProcessTokens() for arguments
9   else if  $t$  is an operator or variable then
10    Add  $t$  to  $\mathcal{P}$ 
11 Step 3: Convert Formula to Prefix Notation
12  $\mathcal{F}_{\text{prefix}} \leftarrow \text{InfixToPrefix}(\mathcal{P})$  // Transform logical
   formula from infix to prefix notation
13 Recursively apply InfixToPrefix() for predicate
   arguments
14 Step 4: Determine Sorts
15  $\mathcal{S}_{\text{sorts}} \leftarrow \text{UnifySort}(\mathcal{F}_{\text{prefix}})$  // Assign sorts for
   variables and predicates
16 Step 5: Format Formula for SMT
17  $\mathcal{F}_{\text{SMT}} \leftarrow \text{Parenthesize } \mathcal{F}_{\text{prefix}}$  according to SMT-LIB
   syntax
18 Step 6: Generate SMT File
19  $\mathcal{S} \leftarrow \text{GenerateSMT}(\mathcal{S}_{\text{sorts}}, \mathcal{F}_{\text{SMT}})$ 
20 Include
   • (declare-sort) statements for sorts.
   • (declare-fun) statements for variables and
     predicates.
   • Negation of  $\mathcal{F}_{\text{SMT}}$ .
   • (check-sat) and (get-model) commands.
return  $\mathcal{S}$  // Return the SMT file for use in
formal solvers
```

B Prompt Examples

Note: Additional in-context examples were removed for brevity and denoted '[...]' in the following prompts.

B.1 End-to-end LLM Prompts

Prompt 1. Classifying with in-context examples (Few-shot)

Logical fallacies are common errors in reasoning that undermine the logic of an argument.

A sentence is logically valid if and only if it is not possible for it to be false.

Algorithm 2: UnifySort for Predicate $A(x, y)$

Input: Predicate $A(x, y)$ with arguments and potential instances

Output: Unified sort for predicate A or an error if sorts are incompatible

```
1 Step 1: Declare the Current Sort
2 Initialize the current sort of  $A$ : (NULL, NULL, Bool)
3 Step 2: Process Each Instance of Predicate  $A$ 
4 foreach instance of predicate  $A$  do
5   Step 2.1: Determine Instance Sorts
6   foreach argument  $x_i$  in the instance do
7     if  $x_i$  is a formula then
8       Set  $\text{sort}(x_i) = \text{Bool}$ 
9     else if  $x_i$  is a variable then
10      Set  $\text{sort}(x_i) = \text{sort}(\text{variable})$  // May be
        NULL
11   Step 2.2: Unify Current Sort with Instance Sort
12   foreach statement sort in current and instance sorts
   do
13     if sorts are not NULL and different then
14       Raise an error: Incompatible sorts
15     else if current sort is NULL and instance sort is
       not NULL then
16       Update current sort:
17       current_sort  $\leftarrow$  instance_sort
18     else if instance sort is NULL and current sort is
       not NULL then
19       Update variable sort to match current sort
19 return Unified sort of predicate  $A$  or error if sorts are
incompatible
```

Here are some examples of classifying sentences as logical fallacies or valid sentences:

Example 1:

Input: "I met a tall man who loved to eat cheese, now I believe all tall people like cheese"

Answer: Logical Fallacy

[...]

Now, classify the following sentence. Answer with either "Logical Fallacy" or "Valid" at the start of your answer.

Input:

1216

B.2 Intermediate NL2FOL Prompts

1217

Prompt 2. Extracting claim and implication

Here are some examples of extracting claims and implications from an input paragraph. There can be multiple claims but only one implication.

Input: "I met a tall man who loved to eat cheese, now I believe all tall people like cheese."

Output:

Claim: "A tall man loves cheese."

Implication: "All tall people like cheese."

[...]

Do not use any subordinating conjunctions in the implication. Replace pronouns with the appropriate nouns so that there are no pronouns. Now extract the claim and implication for the following input.

Input:

1218

Prompt 3. Getting referring expressions

You are given a sentence. Referring expressions are noun phrases, pronouns, and proper names that refer to some individual objects that have some properties associated with them. Here are some examples of finding referring expressions in a sentence:

Input: "A tall man loved cheese"

Referring expressions: A tall man

[...]

Now, find the referring expressions for the following input:

1219

Prompt 4. Getting entity relations

Please determine the relationship between the two entities provided below. Choose the number corresponding to the statement that best describes their relationship:

1. "[Entity A]" is equal to "[Entity B]".

1220

- 2. "[Entity A]" is a subset of "[Entity B]".
- 3. "[Entity B]" is a subset of "[Entity A]".
- 4. "[Entity A]" is not related to "[Entity B]".

Instructions:

- Equality check: If the two entities are equal (case-insensitive after stripping whitespace), select statement 1.
- Subset determination: If they are not equal, assess whether one entity is a subset of the other based on general knowledge and logical reasoning.
 - If "[Entity A]" is a subset of "[Entity B]", select statement 2.
 - If "[Entity B]" is a subset of "[Entity A]", select statement 3.
- Unrelated entities: If none of the above statements accurately describes the relationship.

Here are some examples:

Example 1:

Entity A: "dogs"

Entity B: "animals"

Analysis: All dogs are animals, so "dogs" is a subset of "animals".

Answer: 2

[...]

Entities:

- Entity A:
- Entity B:

Your Task:

- Analyze the relationship between "Entity A" and "Entity B" based on the instructions.
- Provide only the number (1, 2, 3, or 4) that corresponds to the statement you have selected.

1221

Prompt 5. Getting properties (claim)

Given a sentence, and the referring expressions of that sentence. Properties are anything that describes a relationship between two referring expressions, or they may describe a trait of a referring

1222

expression. These properties are essentially predicates in first-order logic.

Here are some examples of finding properties in a sentence:

Example 1:

Input sentence: A tall man loves cheese
Referring expressions: tall man: a, cheese: b
Properties: IsTall(x), LovesCheese(x)

[...]

Now extract the properties for the following input:

sentence into a first-order logical form. Use \rightarrow to represent implies, $\&$ to represent and, \vee to represent or and \neg to represent negations.

Example 1:

Input Sentence: A tall man loves cheese
Referring Expressions: A tall man: x
Properties: IsTall(x), LovesCheese(x)
Logical Form: IsTall(x) $\&$ LovesCheese(x)

[...]

1226

Prompt 6. Getting property relations

You are given two logical clauses. Your task is to identify whether or not the first clause entails the second clause, taking into account external knowledge or 'common sense'. Also, take into account the context from the input sentence.

Here are some examples:

Example 1:

Input sentence: A boy is jumping on skateboard in the middle of a red bridge. Thus, the boy does a skateboarding trick.
Clause 1: JumpsOn(boy,skateboard)
Clause 2: Does(boy, skateboarding_trick)
Answer: ENTAILMENT

[...]

Now given the following clauses. identify whether the first clause entails the second clause.

Prompt 7. Retrieving FOL expression

Given a sentence, the referring expressions of that sentence, and properties which are associated with the referring expressions. Use the given properties to convert the