

# LogicBench: Towards Systematic Evaluation of Logical Reasoning Ability of Large Language Models

Mihir Parmar<sup>1</sup> Nisarg Patel<sup>1</sup> Neeraj Varshney<sup>1</sup> Mutsumi Nakamura<sup>1</sup> Man Luo<sup>1</sup>  
Santosh Mashetty<sup>1</sup> Arindam Mitra<sup>2</sup> Chitta Baral<sup>1</sup>

<sup>1</sup>Arizona State University <sup>2</sup>Microsoft Research

{mparmar3, nppatel17, chitta}@asu.edu

## Abstract

Recently developed large language models (LLMs) have been shown to perform remarkably well on a wide range of language understanding tasks. But, can they really “reason” over the natural language? This question has been receiving significant research attention and many reasoning skills such as common-sense, numerical, and qualitative have been studied. However, the crucial skill pertaining to ‘logical reasoning’ has remained underexplored. Existing work investigating this reasoning ability of LLMs has focused only on a couple of inference rules (such as modus ponens and modus tollens) of propositional and first-order logic. Addressing the above limitation, we comprehensively evaluate the logical reasoning ability of LLMs on 25 different reasoning patterns spanning over propositional, first-order, and non-monotonic logics. To enable systematic evaluation, we introduce *LogicBench*, a natural language question-answering dataset focusing on the use of a single inference rule. We conduct detailed analysis with a range of LLMs such as GPT-4, ChatGPT, Gemini, Llama-2, and Mistral using chain-of-thought prompting. Experimental results show that existing LLMs do not fare well on *LogicBench*; especially, they struggle with instances involving complex reasoning and negations. Furthermore, they sometimes overlook contextual information necessary for reasoning to arrive at the correct conclusion. We believe that our work and findings facilitate future research for evaluating and enhancing the logical reasoning ability of LLMs<sup>1</sup>.

## 1 Introduction

Large language models such as GPT-4, ChatGPT, Google Gemini, Llama-2 (Touvron et al., 2023), and Mistral (Jiang et al., 2023) have made remarkable progress in NLP research enabling machines

<sup>1</sup>Data and code are available at <https://github.com/Mihir3009/LogicBench>

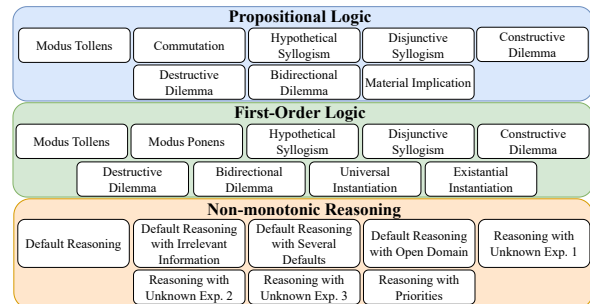


Figure 1: Comprehensive representation of different inference rules and reasoning patterns covered by propositional, first-order, and non-monotonic logics. *Exp.* indicates Expectation

to perform a variety of language tasks that were previously thought to be exclusive to humans (OpenAI, 2023; Brown et al., 2020; Zhao et al., 2023). However, the ability of these LLMs to reason “logically” over natural language text remains under-explored, even though logical reasoning is a fundamental aspect of intelligence and a crucial requirement for many practical applications, such as question-answering systems (Khashabi, 2019) and conversational agents (Beygi et al., 2022). Although several datasets have been proposed (Clark et al., 2021; Tian et al., 2021; Joshi et al., 2020; Saeed et al., 2021) to evaluate the logical reasoning capabilities of LLMs, these datasets are limited in their scope by (1) not evaluating logical reasoning independently of other forms of reasoning such as LogiQA (Liu et al., 2021a) and ReClor (Yu et al., 2020); and (2) evaluating only a single type of logic and covering only few logical inference rules as done in FOLIO (Han et al., 2022) and ProntoQA (Saparov and He, 2023). Thus, our aim in this work is to address the lacuna of having a more comprehensive set of inference rules for evaluating the logical reasoning ability of LLMs.

To this end, we introduce a systematically created question-answering dataset for the evaluation of logical reasoning ability using a single inference

rule, called *LogicBench*. Besides evaluating the logical reasoning ability of LLMs, by evaluating models on single inference rules, we can also gain insights into the frequency of text sequences corresponding to these rules in the pre-training data and their impact on model performance. As illustrated in Figure 1, *LogicBench* includes a total of 25 reasoning patterns across ‘propositional, first-order, and non-monotonic’ logics. To the best of the authors’ knowledge, this is the first work to study non-monotonic reasoning, as well as various inference rules in propositional and first-order logics including hypothetical and disjunctive syllogism; and bidirectional, constructive, and destructive dilemmas in NLP domain. To evaluate LLMs using *LogicBench*, we formulate two different tasks: (i) a Binary Question-Answering (BQA) task in which the context comprises logical statements and the models have to determine whether a conclusion given in the question is logically entailed by the context, and (ii) a Multiple-Choice Questions-Answering (MCQA) task where models must select the most appropriate logical conclusion from four distinct options, based on the provided context. The rationale behind having BQA and MCQA tasks is that it provides systematic standard metric-based evaluation (i.e., direct comparison of LLMs’ performance in terms of accuracy), which could be more challenging with open-ended question-answer formats. Examples instances of various reasoning patterns are presented in Table 4 and App. C.

To construct *LogicBench*, we use a three-stage procedure (refer to §3). In the first stage, we prompt GPT-3.5 to generate a variety of coherent natural language sentences having different ‘ontologies’ (i.e., a collection of concepts as car, person, and animals) and their corresponding negations (refer to §3.2.1). In the second stage, we generate (*context*, *question*) pairs where the context represents a natural language narrative consisting of logical statements, and the question is formulated to exhibit the logical conclusion derived from the context. In the third stage, we generate task-specific data instances (i.e., (*context*, *question*, *answer*) triplets).

We conduct a comprehensive evaluation with a range of LLMs on *LogicBench* including GPT-4, ChatGPT (GPT-3.5-Turbo), Gemini-Pro, Llama-2-7B-Chat, and Mistral-7B-Instruct using chain-of-thought (Wei et al., 2022). In particular, we measure the accuracy of LLMs predictions on both BQA and MCQA tasks. Our experiments result in several interesting findings such as LLMs often

Dataset	Logic Covered			Inference Rules/Axioms Provided with Data	Generation Code Available
	PL	FL	NM		
Ruletaker	✗	✓	✗	✗	Human-annotated
LogicNLI	✗	✓	✗	✗	Semi-automated
ProofWriter	✓	✓	✗	✗	✗
FOLIO	✗	✓	✗	✗	Human-annotated
SimpleLogic	✓	✗	✗	✗	✓
ProntoQA	✗	✓	✗	✓	✓
<b>LogicBench</b>	✓	✓	✓	✓	✓

Table 1: *LogicBench*’s comparison with current datasets

struggle to reason over complex logical contexts that involve only a single inference rule and encounter difficulties with inference rules involving negations. Experimental results reveal that these models struggle with respect to many of the inference rules and patterns, suggesting significant room for improvement in their logical reasoning abilities. To further demonstrate the use of *LogicBench*, we synthetically augment it and fine-tune T5-large. Our preliminary results (App. I) show that this improves the logical reasoning ability of existing models leading to performance improvement on other logic datasets, LogicNLI, and FOLIO (~ 2% on an average), and shows competitive performance on LogiQA and ReClor.

## 2 Related Work

As LLMs continue to evolve rapidly, it becomes increasingly crucial to evaluate their diverse reasoning capabilities, as well as those of forthcoming LLMs. LogiQA (Liu et al., 2021a) and ReClor (Yu et al., 2020) have made notable contributions by compiling multichoice questions from standardized examinations that demand diverse forms of logical reasoning. In contrast to *LogicBench*, these datasets involve mixed forms of reasoning and do not focus on assessing logical reasoning in isolation.

A few past attempts have been made to evaluate only logical reasoning while excluding other forms of reasoning. For example, CLUTTER (Sinha et al., 2019) covers inductive reasoning, (Hahn et al., 2021) covers temporal logic, and Ruletaker (Clark et al., 2021) evaluates whether a transformer-based model emulates deductive reasoning over synthetically generated statements in a limited setting. LogicNLI (Tian et al., 2021) introduced a diagnostic benchmark for FOL reasoning, with the dataset constructed by automatically generating logic expressions and replacing the entity and attribute placeholders.

Our proposed dataset is similar (in terms of task formulation) to ProofWriter (Tafjord et al., 2021),

Names	Propositional Logic	Extension to a (restricted) First-order Logic
MP	$((p \rightarrow q) \wedge p) \vdash q$	$(\forall x(p(x) \rightarrow q(x)) \wedge p(a)) \vdash q(a)$
MT	$((p \rightarrow q) \wedge \neg q) \vdash \neg p$	$(\forall x(p(x) \rightarrow q(x)) \wedge \neg q(a)) \vdash \neg p(a)$
HS	$((p \rightarrow q) \wedge (q \rightarrow r)) \vdash (p \rightarrow r)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (q(x) \rightarrow r(x)))) \vdash (p(a) \rightarrow r(a))$
DS	$((p \vee q) \wedge \neg p) \vdash q$	$(\forall x(p(x) \vee q(x)) \wedge \neg p(a)) \vdash q(a)$
CD	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (p \vee r)) \vdash (q \vee s)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (p(x) \vee r(x)))) \vdash (q(a) \vee s(a))$
DD	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (\neg q \vee \neg s)) \vdash (\neg p \vee \neg r)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (\neg q(x) \vee \neg s(x)))) \vdash (\neg p(a) \vee \neg r(a))$
BD	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (p \vee \neg s)) \vdash (q \vee \neg r)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (p(x) \vee \neg s(x)))) \vdash (q(a) \vee \neg r(a))$
CT	$(p \vee q) \vdash (q \vee p)$	-
MI	$(p \rightarrow q) \vdash (\neg p \vee q)$	-
EG	-	$P(a) \Rightarrow \exists x P(x)$
UI	-	$\forall x A \Rightarrow A\{x \mapsto a\}$

Table 2: Inference rules and (two) axioms that establish the relationship between premises and conclusions. MP: Modus Ponens, MT: Modus Tollens, HS: Hypothetical Syllogism, DS: Disjunctive Syllogism, CD: Constructive Dilemma, DD: Destructive Dilemma, BD: Bidirectional Dilemma, CT: Commutation, MI: Material Implication, EG: Existential Generalization, UI: Universal Instantiation

FOLIO (Han et al., 2022), and ProntoQA (Saparov and He, 2023) which are QA datasets designed to test reasoning ability. ProofWriter provides multi-hop proofs for each example, while FOLIO gives diverse and complex logical expressions, however, it is only limited to FOL. ProntoQA (Saparov and He, 2023) provides explanation and reasoning steps but is limited to modus ponens in FOL. Nevertheless, several crucial attributes motivated us to create *LogicBench* (see Table 1 for comparison). Additional datasets for evaluating logical reasoning also exist such as SimpleLogic (Zhang et al., 2023) provides a class of logical reasoning problems, TaxiNLI (Joshi et al., 2020) introduces logical taxonomy in NLI task and RuleBert (Saeed et al., 2021) covers only soft logical rules. In summary, *LogicBench* evaluates logical reasoning in isolation and provides diverse inference rules and logic types compared to existing datasets. Extended related work is discussed in App. B.

### 3 LogicBench

#### 3.1 Logics Types

**Propositional Logic (PL)** Propositional logic employs a collection of statements or propositions (denoted as  $\mathcal{P} = p_1, p_2, \dots, p_n$ , where  $p_i$  represents a proposition) and builds upon them using logical connectives such as ‘ $\wedge$ ’, ‘ $\vee$ ’, ‘ $\rightarrow$ ’, ‘ $\leftrightarrow$ ’, and ‘ $\neg$ ’. Several inference rules for propositional logic have been defined using which given a set of premises, one can derive a sound conclusion. For example, let’s consider two propositions:  $p_1$ , which states "It is raining," and  $p_2$ , which states "It is cloudy." From these propositions, we can

construct a context/knowledge base (KB) with two premises: (1)  $p_1 \rightarrow p_2$  and (2)  $p_1$ . With this KB, we can conclude  $p_2$ . This inference rule is written as  $((p_1 \rightarrow p_2) \wedge p_1) \vdash p_2$  and is known as ‘Modus Ponens’. In our study, we explore nine distinct inference rules of propositional logic, extensions of seven of them with one-variable and a universal quantifier, and two axioms of first-order logic as shown in Table 2. These inference rules provide a proper framework for deriving valid conclusions.

**First-order Logic (FOL)** In this work, we consider a restricted set of logical axioms for FOL that utilize quantifiers,  $\forall$  (universal quantifier) and  $\exists$  (existential quantifier). The universal quantifier ( $\forall$ ) denotes that a statement holds true for all instances within a specific category. In contrast, the existential quantifier ( $\exists$ ) indicates that a statement is true for at least one instance within its scope. For instance, a simple extension of propositional ‘Modus Ponens’ is an inference rule where given the premises  $\forall(p(x) \rightarrow q(x))$  and  $p(a)$ , we conclude  $q(a)$  (e.g., given “All kings are greedy” and “Sam is a king”, we can conclude “Sam is greedy”). Here, we explore two axioms (EG and UI - in detail in App. C.3) and various inference rules that incorporate the quantifiers (shown in Table 2).

**Non-monotonic (NM) Reasoning** In this work, we analyze a range of logical reasoning templates in NM logics involving “Default Reasoning,” “Reasoning about Unknown Expectations,” and “Reasoning about Priorities.” These templates are inspired by the compilation (Lifschitz, 1989) made in 1989 to evaluate the abilities of various non-

Basic Default Reasoning	Default Reasoning with Irrelevant Information
<p>Context: Block A and block B are both heavy objects that are typically found on the table. However, there is a possibility that block A might not follow this usual convention. It is important to note this exception.</p> <p>Conclusion: B is on the table.</p>	<p>Context: In a room filled with various objects, two heavy blocks, block A and block B, stand out. Normally, heavy blocks like these are placed on the table, but surprisingly, block A is not found on the table. On the other hand, block B grabs attention with its vibrant red color.</p> <p>Conclusion: B is on the table.</p>
Reasoning about Unknown Expectations	Reasoning about Priorities
<p>Context: In this situation, there are three heavy blocks: A, B, and C. Typically, heavy blocks are found on the table. However, it is known that at least one of the blocks, either A or B, is not currently on the table.</p> <p>Conclusion: C is on the table. Exactly one of A, B is not on the table.</p>	<p>Context: John confidently states that the vehicle is situated in the driveway, while Sara adamantly counters, asserting that it is not parked inside the garage.</p> <p>Conclusion: If John’s evidence is more reliable than Sara’s then the car is parked in the driveway.</p>

Table 3: Illustrative examples of non-monotonic reasoning adapted from (Lifschitz, 1989).

monotonic logics that were being developed at that time. Below Table 3 shows examples of NM reasoning. Additional examples are given in App. C.4.

A key aspect of NM logics is to formalize notions such as “normally,” “typically,” and “usually” that are not directly formalizable using classical quantifiers in the first-order setting. The general rule “Heavy blocks are normally located on the table” does not imply that “All heavy blocks are always located on the table”. Rather, this rule allows for exceptions. Our work explores various NM reasoning patterns, as depicted in Figure 1, to delve deeper into the nuances of this type of reasoning.

### 3.2 Data Creation

Our data creation procedure, illustrated in Figure 2, consists of three stages:

**1. Sentence Generation:** Starting with a given prompt, we generate coherent sentences and their negations that incorporate different ontologies.

**2. NL Conversion:** Pairs of (*context*, *question*) are generated using pre-defined templates from which context is then converted to a natural language narrative using the prompt.

**3. Task Instance Generation:** Task-specific (*context*, *question*, *answer*) triplets are generated. BQA requires answers in the form of “yes” or “no”, whereas MCQA involves selecting one correct option from a set of four. We generate semantically preserving and inverting variations of these triplets to add more diversity for BQA.

Examples of generated data corresponding to each logic type and reasoning patterns are presented in App. C.

#### 3.2.1 Sentence Generation

Here, the first step is to generate sentences with diverse *ontologies*. An ontology represents a col-

lection of concepts (e.g. car, person, animals, etc.) along with their corresponding associated properties. To generate these sentences, we prompt the GPT-3.5 model with instructions tailored for each inference rule (more details in App. A).

An example of a prompt corresponding to the ‘Modus Tollens’ from PL is presented in App. A for better illustration. Note that our objective at this stage is not to generate logical sentences but rather to generate a diverse and coherent set of sentences that encompass various concepts. We also create a negation sentence corresponding to each generated sentence<sup>2</sup>. In this work, the scope of generating negations is simple (refer to App. C for examples), however, negations can be more complicated in the case of logic. These generated sentences will be combined with logical connectives in a later stage to form context and questions.

#### 3.2.2 NL Conversion

Here, the NL conversion is accomplished using two steps. First, we leverage the formal expressions of reasoning patterns to create templates that establish the desired NL formulation for each logical connective (i.e., templated context). Second, we prompt GPT-3.5 to transform the templated context into a story/narrative-based context, enhancing its naturalness. For instance, implication: “ $p \rightarrow q$ ” is expressed as “If  $p$ , then  $q$ ”, conjunction: “ $p \wedge q$ ” as “ $p$  and  $q$ ”, and disjunction: “ $p \vee q$ ” as “At least one of the following is true: (1)  $p$  and (2)  $q$ . Note that we do not know which of (1) and (2) is true. It is possible that only (1) is true, or only (2) is true, or both are true.” since understanding the logical implication of ‘or’ when integrated into logical formulations posed challenges to both humans and

<sup>2</sup>We use <https://github.com/dml1s/negate> to generate negated sentences



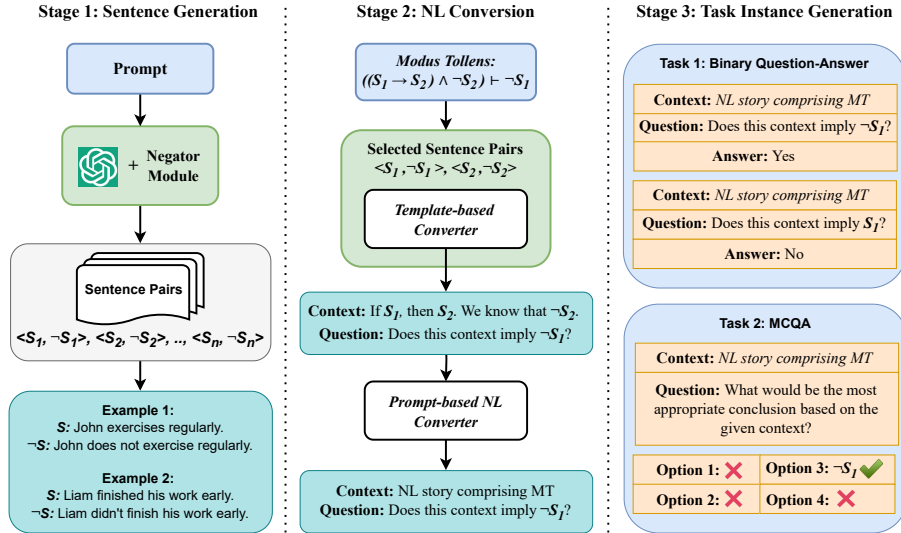


Figure 2: Schematic representation of three-stage procedure for data creation. NL: Natural Language

models. With these established formulations, we proceed to utilize the sentences generated in §3.2.1 to create the templated context and questions corresponding to reasoning patterns.

Then, the templated context is converted into a narrative-based context using a prompt-based converter, enhancing its naturalness. The prompt-based converter (essentially, prompting GPT-3.5) ensures that the context is no longer templated yet follows the logical connection between sentences as mentioned in the logical rule (further details are presented in App. D). For instance, let’s consider the “Modus Tollens” from PL ( $((p \rightarrow q) \wedge \neg q) \vdash \neg p$ ), and the “Bidirectional Dilemma” from FOL ( $(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x))) \wedge (p(a) \vee \neg s(a))) \vdash (q(a) \vee \neg r(a))$ ). For these rules, Table 4 presents examples of logical templated and narrative-based context, question, and task instances for both, BQA and MCQA. App. D showcases further examples corresponding to each inference rule and patterns from *LogicBench*.

### 3.2.3 Task Instance Generation

After generating the context and questions in §3.2.2, we generate *(context, question, answer)* triplets for both tasks: (i) BQA, and (ii) MCQA. Here, narrative-based context is similar for both BQA and MCQA tasks, only the format of *(question, answer)* pairs are different.

**BQA** We generate semantically preserving and inverting variations of questions. Let’s consider the example of “Modus Tollens” from Table 4, having question as: “Does this imply that Liam didn’t

finish his work early?” In this question, we observe one proposition:  $s_1$ , representing the statement “Liam didn’t finish his work early,” can be used to create another question that did not follow the logical rule MT. We can create two possible tuples:  $\langle \neg s_1, yes \rangle, \langle s_1, no \rangle$ . Each tuple has a question-answer combination using proposition  $s_1$ . Moreover, we do not generate variations for the context since it offers no substantial diversity in the dataset. For question variations, we replace the variation of proposition  $\langle s_1 \rangle$  in the original question with the corresponding tuples to add diversity to *LogicBench*. The process allows us to create more variations of the question for BQA tasks, as illustrated in Figure 2 (Step 3 - Task 1).

**MCQA** A prompt-based approach is used to create different incorrect options. For rule MT, as shown in Table 4, option 1 is a correct option that follows the rule MT logically ( $\langle \neg s_1 \rangle$ ), while the other three options are generated using prompting in a way that does not follow the rule being incorrect options. In addition to the options, the question is also replaced by a randomly selected question from the set of five questions. As seen in Figure 2 (Step 3 - Task 2), there is only one correct option out of four given options. More details related to incorrect option generation and instances examples for both tasks are in the App. C, and E.2.

### 3.3 Statistics and Qualitative Analysis

**Statistics** We introduce two versions of our proposed dataset:  $LogicBench(Eval)_{BQA}$  and  $LogicBench(Eval)_{MCQA}$ . Statistics of both ver-

Generated Sentences in Stage 1	Context	Binary QA	MCQA
<p>Inference rule: <b>MT</b></p> <p>p: Liam finished his work early.  <math>\neg</math>p: Liam did not finish his work early.  q: he will order pizza for dinner.  <math>\neg</math>q: he will not order pizza for dinner.</p>	<p><b>Templatized Context:</b> If Liam finishes his work early, then he will order pizza for dinner. He won't order pizza for dinner.</p> <p><b>NL Context:</b> Liam knows that if he finishes his work early for the day, he will order pizza for dinner. However, on this particular day, he decided against ordering pizza.</p>	<p><b>Question 1:</b> Does this imply that Liam didn't finish his work early? (Yes)</p> <p><b>Question 2:</b> Does this imply that Liam finishes his work early? (No)</p>	<p><b>Question:</b> Based on the context, what conclusion would be deemed most suitable?</p> <p><b>Option 1:</b> Liam didn't finish his work early. (Yes)</p> <p><b>Option 2:</b> Sarah had already ordered Chinese takeout. (No)</p> <p><b>Option 3:</b> Rebecca finished her work early. (No)</p> <p><b>Option 4:</b> Liam decided to order sushi instead. (No)</p>
<p>Inference rule: <b>BD</b></p> <p>p(x): someone drinks lots of water  q(x): they will feel hydrated  r(x): they eat too much sugar  s(x): they will experience a sugar crash  p(a): Jane drinks lots of water  <math>\neg</math>p(a): Jane does not drink lots of water  q(a): she will feel hydrated  <math>\neg</math>q(a): she will not feel hydrated  r(a): she eats too much sugar  <math>\neg</math>r(a): she does not eat too much sugar  s(a): she will experience a sugar crash  <math>\neg</math>s(a): she will not experience a sugar crash</p>	<p><b>Templatized Context:</b> If someone drinks lots of water, then they will feel hydrated. If they eat too much sugar, then they will experience a sugar crash. We know that at least one of the following is true (1) Jane drinks lots of water and (2) she won't experience a sugar crash. Note that we do not know which ones of (1) and (2) are true. It might be the case that only (1) is true, or only (2) is true or both are true.</p> <p><b>NL Context:</b> If someone consumes a significant amount of water, they will experience a state of hydration. Conversely, if excessive amounts of sugar are ingested by them, a sugar crash will ensue. It is known that at least one of the following statements is true: either the Jane consumes ample water or she will not experience a sugar crash. However, the actual veracity of either statement remains ambiguous, as it could be the case that only the first statement is true, only the second statement is true, or both statements are true.</p>	<p><b>Question 1:</b> Can we say at least one of the following must always be true? (a) she will feel hydrated and (b) she doesn't eat too much sugar (Yes)</p> <p><b>Question 2:</b> Can we say at least one of the following must always be true? (a) she won't feel hydrated and (b) she eats too much sugar (No)</p> <p><b>Question 3:</b> Can we say at least one of the following must always be true? (a) she will feel hydrated and (b) she eats too much sugar (No)</p> <p><b>Question 4:</b> Can we say at least one of the following must always be true? (a) she won't feel hydrated and (b) she doesn't eat too much sugar (No)</p>	<p><b>Question:</b> Taking into account the context provided, what conclusion would be the most appropriate?</p> <p><b>Option 1:</b> If Jane consumes ample water, she will experience a sugar crash (No)</p> <p><b>Option 2:</b> John will feel hydrated or he won't experience a sugar crash (No)</p> <p><b>Option 3:</b> Jane will feel hydrated or she doesn't eat too much sugar (Yes)</p> <p><b>Option 4:</b> Jane won't feel hydrated or she will eat too much sugar (No)</p>

Table 4: Illustrative examples of logical context and questions created using sentences that are generated in the first stage (§3.2.1) for both tasks, i.e., BQA and MCQA.

sions are presented in Table 5. For  $LogicBench(Eval)_{BQA}$ , out of 1520, 520 samples are for ‘yes’ and 1000 samples are for ‘no’ labels. For  $LogicBench(Eval)_{MCQA}$ , there are 20 unique samples present for each rule, thus in total 500 unique samples. Furthermore, we synthetically augmented  $LogicBench(Eval)$  for training purposes (i.e.,  $LogicBench(Aug)$ ) which consists of 150 unique data samples for each rule for BQA, resulting in a total of 12908 data samples including variations.

Dataset	# of Instances per Axiom	Total # of Instances	Total # of Instances (Including Variations)
$LogicBench(Eval)_{BQA}$	20	500	1520
$LogicBench(Eval)_{MCQA}$	20	500	500

Table 5: Statistics of  $LogicBench(Eval)$

**Data Validation** Throughout the data generation phase of  $LogicBench(Eval)$ , the authors conduct a review of the logical formations to ensure they follow the intended logical structure. We examine each narrative for any potential discrepancies, ensuring that they are logically sound and correctly represent the intended relationships between propositions. In addition to the logical formation, we also dedicated considerable effort to eliminating typos and validating the grammar. We also analyze the diversity in terms of different ontology and the logical nature of the  $LogicBench(Eval)$  (presented in App. C.1). We mitigate errors encountered during the validation step (presented in App. F).

## 4 Results and Analysis

### 4.1 Experimental Setup

**Task Formulation** For BQA, let us consider a set of data instances  $\mathcal{I}_{r,L}$  corresponding to the inference rule  $r$  and logic type  $L$ . In this set,  $i^{th}$  instance

is represented as  $\mathcal{I}_{r,L}^i = \{(c_i, Q_i)\}$  where  $c_i$  represents narrative context and  $Q_i = \{q_1, q_2, \dots, q_n\}$  represents set of question and its variations corresponding to  $i^{th}$  instance. As discussed in §3, each context ( $c$ ) represents logical rules (e.g., All cats have fur. Tom is a cat.) and question ( $q$ ) represents the conclusion (e.g., Does Tom have fur?). To each context and question pair, i.e.,  $\langle c, q \rangle$ , we assign a label from the set  $\mathcal{Y} = \{Yes, No\}$ . We assign a label *Yes* if the conclusion logically entails the context, otherwise, assign a label *No*. To evaluate any LLMs on this setup, we provide  $\langle p, c, q \rangle$  as input to predict a label from  $\mathcal{Y}$  where  $p$  is a natural language prompt. In the set  $\mathcal{I}_{r,L}$  for MCQA,  $i^{th}$  instance is represented as  $\mathcal{I}_{r,L}^i = \{(c_i, q_i, O_i)\}$  where  $c_i$  represents narrative context and  $q_i$  represents question and  $O_i = \{o_1, o_2, o_3, o_4\}$  represents four option choices. To each context and question pair, i.e.,  $\langle c, q \rangle$ , we assign a label from the set  $\mathcal{Y} = \{o_1, o_2, o_3, o_4\}$ . We assign a label  $o_1$  if the correct conclusion is presented in the first option, and likewise for other labels. To evaluate any LLMs on this setup, we provide  $\langle p, c, q, o \rangle$  as input to predict a label from  $\mathcal{Y}$ .

**Experiments** We evaluate a range of prompting models including GPT-4, ChatGPT (GPT-3.5-Turbo), Google Gemini-Pro, Llama-2-7B-Chat, and Mistral-7B-Instruct-v0.2. Each model is evaluated in a zero-shot setting where the chain-of-thought prompt is provided to the model without any in-context examples. This approach allows us to determine LLM’s inherent ability to do logical reasoning (based on pre-training), as we can not expect that various logical inference rules/patterns will always be made part of prompts. However, we

Type	Rules	Llama-2		Mistral		Gemini		ChatGPT		GPT-4	
		$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$
PL	HS	100 <sub>0.00</sub>	47.35 <sub>0.03</sub>	98.63 <sub>0.01</sub>	64.24 <sub>0.07</sub>	99.31 <sub>0.01</sub>	59.60 <sub>0.01</sub>	100 <sub>0.00</sub>	66.82 <sub>0.04</sub>	100 <sub>0.00</sub>	84.75 <sub>0.06</sub>
	DS	44.81 <sub>0.05</sub>	56.82 <sub>0.02</sub>	50.77 <sub>0.03</sub>	77.66 <sub>0.05</sub>	68.21 <sub>0.04</sub>	91.65 <sub>0.04</sub>	51.26 <sub>0.04</sub>	80.71 <sub>0.09</sub>	79.96 <sub>0.06</sub>	100 <sub>0.00</sub>
	CD	79.94 <sub>0.03</sub>	25.47 <sub>0.01</sub>	83.37 <sub>0.03</sub>	72.30 <sub>0.12</sub>	85.56 <sub>0.04</sub>	28.15 <sub>0.01</sub>	90.75 <sub>0.06</sub>	38.63 <sub>0.04</sub>	92.85 <sub>0.01</sub>	66.12 <sub>0.05</sub>
	DD	82.22 <sub>0.16</sub>	25.22 <sub>0.01</sub>	71.26 <sub>0.02</sub>	14.16 <sub>0.03</sub>	75.48 <sub>0.04</sub>	25.22 <sub>0.02</sub>	71.28 <sub>0.05</sub>	23.47 <sub>0.03</sub>	84.51 <sub>0.02</sub>	42.40 <sub>0.04</sub>
	BD	48.89 <sub>0.43</sub>	25.13 <sub>0.01</sub>	79.68 <sub>0.01</sub>	47.36 <sub>0.11</sub>	86.81 <sub>0.02</sub>	29.51 <sub>0.01</sub>	84.91 <sub>0.05</sub>	33.47 <sub>0.05</sub>	87.86 <sub>0.04</sub>	59.06 <sub>0.13</sub>
	MT	70.54 <sub>0.03</sub>	71.46 <sub>0.05</sub>	47.96 <sub>0.08</sub>	44.96 <sub>0.10</sub>	75.99 <sub>0.05</sub>	81.47 <sub>0.03</sub>	55.63 <sub>0.02</sub>	66.11 <sub>0.07</sub>	55.28 <sub>0.04</sub>	59.05 <sub>0.07</sub>
	MI	78.57 <sub>0.26</sub>	25.34 <sub>0.01</sub>	75.36 <sub>0.01</sub>	25.51 <sub>0.04</sub>	74.69 <sub>0.06</sub>	24.84 <sub>0.04</sub>	81.60 <sub>0.02</sub>	31.79 <sub>0.03</sub>	91.84 <sub>0.01</sub>	39.72 <sub>0.03</sub>
	CT	70.00 <sub>0.12</sub>	24.98 <sub>0.02</sub>	87.38 <sub>0.01</sub>	71.99 <sub>0.03</sub>	88.31 <sub>0.08</sub>	35.11 <sub>0.07</sub>	89.88 <sub>0.05</sub>	43.33 <sub>0.03</sub>	98.59 <sub>0.01</sub>	60.71 <sub>0.06</sub>
	<b>Avg</b>	<b>71.87<sub>0.13</sub></b>	<b>37.72<sub>0.02</sub></b>	<b>74.30<sub>0.03</sub></b>	<b>52.27<sub>0.07</sub></b>	<b>81.79<sub>0.04</sub></b>	<b>46.94<sub>0.03</sub></b>	<b>78.16<sub>0.04</sub></b>	<b>48.04<sub>0.05</sub></b>	<b>86.36<sub>0.03</sub></b>	<b>63.98<sub>0.05</sub></b>
FOL	EG	100 <sub>0.00</sub>	71.67 <sub>0.03</sub>	100 <sub>0.00</sub>	100 <sub>0.00</sub>	85.68 <sub>0.06</sub>	97.97 <sub>0.04</sub>	96.74 <sub>0.03</sub>	96.75 <sub>0.03</sub>	98.41 <sub>0.03</sub>	100 <sub>0.00</sub>
	UI	80.24 <sub>0.05</sub>	56.04 <sub>0.02</sub>	85.24 <sub>0.04</sub>	85.24 <sub>0.04</sub>	90.69 <sub>0.08</sub>	90.29 <sub>0.04</sub>	85.31 <sub>0.02</sub>	96.19 <sub>0.03</sub>	90.16 <sub>0.002</sub>	91.58 <sub>0.03</sub>
	MP	97.44 <sub>0.04</sub>	73.01 <sub>0.04</sub>	96.28 <sub>0.03</sub>	95.24 <sub>0.05</sub>	98.25 <sub>0.03</sub>	96.97 <sub>0.05</sub>	100 <sub>0.00</sub>	91.03 <sub>0.04</sub>	94.74 <sub>0.00</sub>	100 <sub>0.00</sub>
	HS	100 <sub>0.00</sub>	36.83 <sub>0.01</sub>	98.02 <sub>0.00</sub>	65.00 <sub>0.05</sub>	97.70 <sub>0.00</sub>	52.46 <sub>0.04</sub>	95.14 <sub>0.00</sub>	46.74 <sub>0.04</sub>	89.03 <sub>0.01</sub>	61.21 <sub>0.01</sub>
	DS	58.33 <sub>0.14</sub>	52.08 <sub>0.04</sub>	61.69 <sub>0.03</sub>	92.59 <sub>0.13</sub>	75.96 <sub>0.07</sub>	85.82 <sub>0.07</sub>	72.46 <sub>0.10</sub>	94.71 <sub>0.05</sub>	63.34 <sub>0.04</sub>	100 <sub>0.00</sub>
	CD	100 <sub>0.00</sub>	26.36 <sub>0.01</sub>	82.21 <sub>0.02</sub>	82.02 <sub>0.09</sub>	94.71 <sub>0.02</sub>	33.62 <sub>0.02</sub>	87.87 <sub>0.06</sub>	39.37 <sub>0.06</sub>	92.65 <sub>0.01</sub>	74.72 <sub>0.05</sub>
	DD	54.37 <sub>0.19</sub>	23.36 <sub>0.01</sub>	73.78 <sub>0.01</sub>	6.67 <sub>0.11</sub>	82.70 <sub>0.02</sub>	28.96 <sub>0.01</sub>	72.70 <sub>0.08</sub>	26.42 <sub>0.11</sub>	85.54 <sub>0.04</sub>	54.55 <sub>0.08</sub>
	BD	92.31 <sub>0.13</sub>	25.78 <sub>0.01</sub>	75.76 <sub>0.01</sub>	44.44 <sub>0.10</sub>	82.82 <sub>0.08</sub>	28.93 <sub>0.05</sub>	85.50 <sub>0.03</sub>	38.82 <sub>0.02</sub>	84.22 <sub>0.02</sub>	66.67 <sub>0.10</sub>
MT	77.11 <sub>0.05</sub>	88.41 <sub>0.09</sub>	69.01 <sub>0.06</sub>	91.71 <sub>0.01</sub>	71.18 <sub>0.10</sub>	89.58 <sub>0.04</sub>	62.46 <sub>0.03</sub>	92.80 <sub>0.06</sub>	64.69 <sub>0.04</sub>	100 <sub>0.00</sub>	
	<b>Avg</b>	<b>84.42<sub>0.07</sub></b>	<b>50.39<sub>0.03</sub></b>	<b>82.44<sub>0.02</sub></b>	<b>73.66<sub>0.06</sub></b>	<b>86.63<sub>0.05</sub></b>	<b>67.17<sub>0.04</sub></b>	<b>84.24<sub>0.04</sub></b>	<b>69.20<sub>0.05</sub></b>	<b>84.75<sub>0.02</sub></b>	<b>83.19<sub>0.02</sub></b>
NM	DRI	29.06 <sub>0.05</sub>	40.98 <sub>0.04</sub>	61.51 <sub>0.05</sub>	72.22 <sub>0.05</sub>	65.45 <sub>0.03</sub>	78.50 <sub>0.02</sub>	55.77 <sub>0.00</sub>	88.89 <sub>0.09</sub>	81.29 <sub>0.05</sub>	100 <sub>0.00</sub>
	DRS	66.56 <sub>0.03</sub>	18.40 <sub>0.01</sub>	70.13 <sub>0.01</sub>	6.20 <sub>0.01</sub>	68.05 <sub>0.01</sub>	7.20 <sub>0.06</sub>	69.06 <sub>0.01</sub>	0.00 <sub>0.00</sub>	77.59 <sub>0.02</sub>	39.32 <sub>0.06</sub>
	DRD	49.21 <sub>0.06</sub>	50.47 <sub>0.04</sub>	77.64 <sub>0.02</sub>	97.78 <sub>0.04</sub>	68.10 <sub>0.05</sub>	97.43 <sub>0.04</sub>	55.11 <sub>0.02</sub>	100 <sub>0.00</sub>	86.96 <sub>0.00</sub>	100 <sub>0.00</sub>
	DRO	53.25 <sub>0.01</sub>	53.60 <sub>0.02</sub>	68.30 <sub>0.04</sub>	89.56 <sub>0.11</sub>	57.85 <sub>0.04</sub>	74.03 <sub>0.08</sub>	50.85 <sub>0.01</sub>	66.67 <sub>0.58</sub>	55.56 <sub>0.00</sub>	100 <sub>0.00</sub>
	RE1	94.07 <sub>0.06</sub>	27.55 <sub>0.01</sub>	78.79 <sub>0.02</sub>	44.39 <sub>0.03</sub>	84.04 <sub>0.04</sub>	41.98 <sub>0.12</sub>	75.82 <sub>0.00</sub>	31.74 <sub>0.03</sub>	85.44 <sub>0.02</sub>	100 <sub>0.00</sub>
	RE2	54.17 <sub>0.47</sub>	53.62 <sub>0.05</sub>	81.31 <sub>0.04</sub>	85.67 <sub>0.03</sub>	58.48 <sub>0.03</sub>	82.22 <sub>0.17</sub>	64.67 <sub>0.11</sub>	60.68 <sub>0.02</sub>	59.42 <sub>0.00</sub>	100 <sub>0.00</sub>
	RE3	39.05 <sub>0.03</sub>	38.05 <sub>0.10</sub>	78.61 <sub>0.04</sub>	83.59 <sub>0.01</sub>	64.67 <sub>0.08</sub>	79.68 <sub>0.02</sub>	59.05 <sub>0.04</sub>	83.57 <sub>0.08</sub>	83.12 <sub>0.02</sub>	89.08 <sub>0.01</sub>
	RAP	76.82 <sub>0.05</sub>	67.22 <sub>0.03</sub>	51.73 <sub>0.02</sub>	60.48 <sub>0.11</sub>	75.58 <sub>0.04</sub>	96.67 <sub>0.06</sub>	61.44 <sub>0.02</sub>	92.16 <sub>0.03</sub>	66.07 <sub>0.05</sub>	98.72 <sub>0.02</sub>
	<b>Avg</b>	<b>57.78<sub>0.09</sub></b>	<b>43.74<sub>0.04</sub></b>	<b>71.00<sub>0.03</sub></b>	<b>67.49<sub>0.05</sub></b>	<b>67.78<sub>0.04</sub></b>	<b>69.71<sub>0.07</sub></b>	<b>61.47<sub>0.03</sub></b>	<b>65.46<sub>0.10</sub></b>	<b>74.43<sub>0.02</sub></b>	<b>84.75<sub>0.02</sub></b>

Table 6: Evaluation of LLMs in terms of label-wise accuracy on  $LogicBench(Eval)_{BQA}$ , where  $A(Yes)$  and  $A(No)$  denote the accuracy for the  $Yes$  and  $No$  labels, respectively. DRI: Default Reasoning with Irrelevant Information, DRS: Default Reasoning with Several Defaults, DRD: Default Reasoning with a Disabled Default, DRO: Default Reasoning in an Open Domain, RE1: Reasoning about Unknown Expectations I, RE2: Reasoning about Unknown Expectations II, RE3: Reasoning about Unknown Expectations III, RAP: Reasoning about Priorities

do evaluate these models in a few-shot setting, and present the results in App. G.

**Metrics** Here, we evaluate performance in terms of accuracy for both tasks, BQA and MCQA. For the BQA, we measure accuracy corresponding to each label, i.e.,  $A(Yes)$  and  $A(No)$ . We evaluate each model on three different chain-of-thought prompts and report average results across these prompts. All prompts used for experiments are described in App. H.

## 4.2 Main Results

Table 6 and Table 7 represent inference rule-wise performance corresponding to each LLMs for the BCQ and MCQA tasks, respectively. Specifically, Table 6 provides label-wise accuracy ( $A(Yes)$  and  $A(No)$ ) for the BQA task, and Table 7 provides overall accuracy for the MCQA task. Both tables provide valuable insights into the performance of different models on various logic types and lead to several interesting findings. From Table 6, we

can observe that ChatGPT achieves 48.04%, and GPT-4 shows a performance of 63.98%  $A(Yes)$  on average which indicates the challenge of classical logical reasoning (PL) even for larger LLMs such as ChatGPT and GPT-4. Furthermore, we can observe that models struggle more with inference rules of PL compared to FOL and NM. In addition, it is noticeable that each model performs relatively better on questions with a negative response (i.e.,  $No$ ) compared to questions with a positive response (i.e.,  $Yes$ ). This observation suggests that the models struggle to fully comprehend the logical relationship between the context and the conclusion (i.e., lower  $A(Yes)$ ). However, they demonstrate a relatively stronger understanding when the relationship is contradictory in nature (i.e., higher  $A(No)$ ). From Table 7, we can observe that larger models exhibit superior performance in selecting the correct choice to arrive at a logical conclusion. Interestingly, the performance of models decreases when the inference rules are longer or include nega-

Type	Rules	Llama-2	Mistral	Gemini	ChatGPT	GPT-4
PL	HS	86.67 <sub>0.08</sub>	93.33 <sub>0.06</sub>	100 <sub>0.00</sub>	91.67 <sub>0.03</sub>	100 <sub>0.00</sub>
	DS	63.33 <sub>0.12</sub>	60.00 <sub>0.09</sub>	86.67 <sub>0.10</sub>	96.67 <sub>0.06</sub>	95.00 <sub>0.00</sub>
	CD	80.00 <sub>0.05</sub>	70.00 <sub>0.13</sub>	96.67 <sub>0.03</sub>	90.00 <sub>0.00</sub>	100 <sub>0.00</sub>
	DD	43.33 <sub>0.03</sub>	30.00 <sub>0.10</sub>	90.00 <sub>0.05</sub>	73.33 <sub>0.06</sub>	88.33 <sub>0.03</sub>
	BD	51.67 <sub>0.03</sub>	53.33 <sub>0.03</sub>	86.67 <sub>0.03</sub>	68.33 <sub>0.06</sub>	83.33 <sub>0.03</sub>
	MT	31.67 <sub>0.03</sub>	60.00 <sub>0.05</sub>	78.33 <sub>0.08</sub>	73.33 <sub>0.08</sub>	76.67 <sub>0.03</sub>
	MI	33.33 <sub>0.13</sub>	35.00 <sub>0.05</sub>	71.67 <sub>0.06</sub>	63.33 <sub>0.08</sub>	73.33 <sub>0.14</sub>
	CT	73.33 <sub>0.10</sub>	68.33 <sub>0.08</sub>	100 <sub>0.00</sub>	98.33 <sub>0.03</sub>	100 <sub>0.00</sub>
	<b>Avg</b>	<b>57.92<sub>0.07</sub></b>	<b>58.75<sub>0.07</sub></b>	<b>88.75<sub>0.04</sub></b>	<b>81.88<sub>0.05</sub></b>	<b>89.58<sub>0.03</sub></b>
FOL	EI	80.00 <sub>0.00</sub>	85.00 <sub>0.05</sub>	95.00 <sub>0.00</sub>	93.33 <sub>0.03</sub>	100 <sub>0.00</sub>
	UI	63.33 <sub>0.03</sub>	75.00 <sub>0.05</sub>	98.33 <sub>0.03</sub>	91.67 <sub>0.03</sub>	98.33 <sub>0.03</sub>
	MP	85.00 <sub>0.09</sub>	98.33 <sub>0.03</sub>	100 <sub>0.00</sub>	100 <sub>0.00</sub>	100 <sub>0.00</sub>
	HS	61.67 <sub>0.06</sub>	70.00 <sub>0.05</sub>	81.67 <sub>0.03</sub>	73.33 <sub>0.06</sub>	76.67 <sub>0.03</sub>
	DS	43.33 <sub>0.06</sub>	36.67 <sub>0.03</sub>	70.00 <sub>0.05</sub>	78.33 <sub>0.10</sub>	95.00 <sub>0.05</sub>
	CD	75.00 <sub>0.05</sub>	61.67 <sub>0.06</sub>	93.33 <sub>0.03</sub>	80.00 <sub>0.05</sub>	91.67 <sub>0.03</sub>
	DD	36.67 <sub>0.06</sub>	46.67 <sub>0.03</sub>	85.00 <sub>0.05</sub>	71.67 <sub>0.06</sub>	93.33 <sub>0.03</sub>
	BD	35.00 <sub>0.05</sub>	43.33 <sub>0.06</sub>	78.33 <sub>0.06</sub>	66.67 <sub>0.13</sub>	91.67 <sub>0.06</sub>
	<b>Avg</b>	<b>57.96<sub>0.05</sub></b>	<b>64.81<sub>0.04</sub></b>	<b>87.04<sub>0.03</sub></b>	<b>82.41<sub>0.06</sub></b>	<b>91.51<sub>0.04</sub></b>
NM	DRI	38.33 <sub>0.03</sub>	28.33 <sub>0.06</sub>	58.33 <sub>0.08</sub>	66.67 <sub>0.06</sub>	90.00 <sub>0.00</sub>
	DRS	41.67 <sub>0.08</sub>	16.67 <sub>0.10</sub>	45.00 <sub>0.10</sub>	41.67 <sub>0.10</sub>	55.00 <sub>0.05</sub>
	DRD	55.00 <sub>0.00</sub>	50.00 <sub>0.05</sub>	48.33 <sub>0.03</sub>	71.67 <sub>0.10</sub>	80.00 <sub>0.05</sub>
	DRO	21.67 <sub>0.03</sub>	21.67 <sub>0.03</sub>	53.33 <sub>0.03</sub>	38.33 <sub>0.08</sub>	45.00 <sub>0.00</sub>
	RE1	51.67 <sub>0.03</sub>	31.67 <sub>0.08</sub>	70.00 <sub>0.00</sub>	65.00 <sub>0.05</sub>	95.00 <sub>0.05</sub>
	RE2	65.00 <sub>0.05</sub>	75.00 <sub>0.00</sub>	68.33 <sub>0.06</sub>	61.67 <sub>0.03</sub>	66.67 <sub>0.06</sub>
	RE3	31.67 <sub>0.03</sub>	33.33 <sub>0.03</sub>	61.67 <sub>0.08</sub>	70.00 <sub>0.05</sub>	68.33 <sub>0.03</sub>
	RAP	46.67 <sub>0.08</sub>	35.00 <sub>0.09</sub>	33.33 <sub>0.03</sub>	55.00 <sub>0.05</sub>	51.67 <sub>0.03</sub>
	<b>Avg</b>	<b>43.96<sub>0.04</sub></b>	<b>36.46<sub>0.05</sub></b>	<b>54.79<sub>0.05</sub></b>	<b>58.75<sub>0.07</sub></b>	<b>68.96<sub>0.03</sub></b>

Table 7: Evaluation of LLMs in terms of accuracy on *LogicBench(Eval)<sub>MCQA</sub>*.

tions. In contrast to Table 6, for MCQA, LLMs show superior performance for PL and FOL compared to NM. To further investigate these findings and provide a detailed analysis, we perform a thorough study of reasoning chains generated by LLMs and present our insights in the subsequent section.

### 4.3 Analysis and Discussion

**Human Performance** We conduct a human evaluation on a subset of *LogicBench(Eval)* for both tasks, BQA and MCQA. Specifically, we selected 50 unique instances covering all 25 reasoning patterns from *LogicBench(Eval)*. This selection resulted in total instances of 153 <context, question> pairs for BQA, and 50 <context, question, choices> pairs for MCQA. We hired three graduate student volunteers to provide the evaluations. The task instructions given to all three annotators closely resemble the prompts provided to models (App. H). Each instance pair is answered/annotated by three different annotators with 0.785 inter-annotator agreement (measured with raw/observed agreement) for BQA and 0.813 for MCQA.

From the results (Table 8) for BQA, we see that humans achieve more than  $\sim 85\%$  accuracy on various logic types on *LogicBench(Eval)* which indicates the capability of humans to comprehend single-step logical reasoning effectively. From Ta-

Logic Type	BQA		MCQA
	A(No)	A(Yes)	Accuracy
PL	85.42%	84.17%	100%
FOL	90.74%	91.18%	97.03%
NM	72.22%	88.46%	62.05%
<b>Avg.</b>	<b>82.79%</b>	<b>87.94%</b>	<b>86.36%</b>

Table 8: Human performance on three logic types averaged across three annotators for both tasks.

ble 6, we observe that the average performance of all models is below human performance indicating room for improvement in their reasoning capabilities. From Table 7, we make similar observations for MCQA. However, we can see that the performance of NM for the MCQA task remains a challenge for both humans and LLMs.

**Lower performance of LLMs on PL as compared to NM and FOL for BQA.** In the development of AI, NM logic was partly developed to formalize natural language constructs, such as “normally birds fly”, that were not formalizable in a straightforward manner using classical mathematical logics. Thus, while it was difficult for researchers to come up with non-monotonic logics and formalize non-monotonic reasoning, the fact that they were usually motivated by natural language examples, suggests that many of the non-monotonic reasoning aspects are present in the NL text in the wild that is used in the pre-training of the ultra-large LLMs such as GPT4. While from human experience and complexity theory, FOL is harder than PL in general; in the LLM context, the crucial factor becomes what kind of logical sentences LLMs are pre-trained on. It seems that LLMs are pre-trained more on simple FOL sentences than on simple PL sentences (see Appendix I for further discussion). On the other hand, some PL features are perhaps less prevalent in human writing (on which LLMs are pre-trained) - such as Modes Tollens. Table 6 shows that GPT-4 achieves  $\sim 85\%$  accuracy ( $A(Yes)$ ) for simple inference rules such as HS(PL). However, GPT-4 performance dropped to  $\sim 59\%$   $A(Yes)$  for PL(MT).

**Negations are hard to understand when embedded with logical rules.** Regarding PL and FOL, it is apparent that the models struggle more with the DD and MT inference rules. A closer look at Table 2 reveals that all of these rules include examples where the models need to draw conclusions based



on negated premises. This indicates that the models encounter difficulties when negated premises are introduced. We also analyze the effect of negations on the reasoning chain (see App. I).

**Longer inference rules are still challenging.** Table 6 indicates that the models face challenges when handling longer rules, such as BD, CD, and DD, both in PL and FOL. Hence, it can be concluded that these models struggle with longer logical dependencies in the premise, particularly when a higher number of propositions are present. In the case of NM reasoning, the models exhibit lower performance in DRS, indicating that a higher number of premises often leads to more frequent mistakes.

**LLMs sometimes overlook contextual information.** We investigate the LLMs’ logical reasoning ability in natural language, not in artificial logical formulations. Hence, we note that LLMs sometimes hallucinate information and overlook contextual information, leading to incorrect conclusions. To analyze this, we manually examine all the reasoning chains generated for instances sharing the same contexts in both BQA and MCQA tasks. We observe that, although this pattern is not dominant, it affects BQA more than MCQA. For a more in-depth analysis, please refer to App. I.

**Large models are better logical reasoners.** We analyze the results of both smaller (Llama-2-7B and Mistral-7B) and larger (ChatGPT, GPT-4, and Gemini) models. Table 6 and Table 7 show that larger models tend to exhibit higher performance across different types of logic. We further investigate an additional model with an intermediate size: Yi-34B-chat (Young et al., 2024) (results are presented in App. I). When compared to the Llama-2-7B, the Yi-34B model (5x larger than Llama-7B) shows improvement in average performance across three logic types. Similarly, GPT-4 outperforms Yi-34B. This suggests that increasing the model size leads to substantial gains in performance, indicating the influence of larger model capacities on carrying out better logical reasoning.

**Performance of BQA vs. MCQA** From Table 6 and Table 7, we can see the overall performance of LLMs is higher on  $PL_{MCQA}$  compared to  $PL_{BQA}$ . Conversely, the performance is lower on  $NM_{MCQA}$  compared to  $NM_{BQA}$ . For PL, the performance gaps between the CT and DD inference rules primarily contributed to this trend, and

DRO, RAP, and DRD for NM. We analyze the reasoning chains associated with these inference rules and presented our detailed observations in App. I.

**Effect on other logic datasets** We trained the T5-large model on the *LogicBench(Aug)* resulting in a model named LogicT5. Furthermore, we performed fine-tuning on four other logical reasoning datasets: LogiQA, Reclor, LogicNLI, and FOLIO. Further discussion is presented in App. I.

## 5 Conclusions

In this work, we evaluated the logical reasoning ability of LLMs on 25 distinct inference rules and reasoning patterns covering PL, FOL, and NM logics. To this end, we introduced *LogicBench*, a natural language question-answering dataset focusing on evaluating a single inference rule. We devised two tasks using *LogicBench*: (i) BQA, and (ii) MCQA. We evaluated a range of LLMs including GPT-4, ChatGPT, Gemini-Pro, Llama-2, and Mistral on both tasks. Experimental results showed that LLMs do not perform well on *LogicBench*, even though they require the application of only a single inference rule. Furthermore, we also augmented *LogicBench* to *LogicBench(Aug)*, which can be utilized for training purposes. Using *LogicBench(Aug)*, we demonstrated that LLMs trained using it showcase an improved understanding of logical reasoning, resulting in a better performance on existing logic datasets.

## Limitations

While *LogicBench* encompasses 25 distinct inference rules spanning three logic types (significantly more than any previous study) to comprehensively evaluate the logical reasoning capabilities of LLMs, it can be further extended by incorporating additional inference rules and logic types. However, with respect to first-order logic and logics with quantified variables, there can be an infinite number of such rules. In this study, we focused solely on evaluating model performance using a single inference rule; however, an interesting future direction can be enhancing the depth of reasoning complexity (i.e., multi-step reasoning) by incorporating combinations of inference rules to derive conclusions. We also note that this research is limited to the English language and can be extended to multilingual scenarios for evaluating the logical reasoning ability of LLMs.

## Ethics Statement

We have used AI assistants (Grammarly and ChatGPT) to address the grammatical errors and rephrase the sentences.

## Acknowledgement

We thank the anonymous reviewers for constructive suggestions, and the computer science graduate students of Arizona State University (ASU) who helped with the human annotations. We extend our gratitude to the Research Computing (RC) at ASU for providing computing resources for experiments. We acknowledge support by a 2023 Spring Amazon Research Award (ARA), and an award by Cisco via Silicon Valley Foundation.

## References

- Jean-Philippe Bernardy and Stergios Chatzikyriakidis. 2020. [Fracas: Temporal analysis](#).
- Sajjad Beygi, Maryam Fazel-Zarandi, Alessandra Cervone, Prakash Krishnan, and Siddhartha Jonnalagadda. 2022. [Logical reasoning for task oriented dialogue systems](#). In *Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)*, pages 68–79, Dublin, Ireland. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Xiang Chen, Xin Xie, Ningyu Zhang, Jiahuan Yan, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021. [Adaprompt: Adaptive prompt-based finetuning for relation extraction](#). *arXiv e-prints*, pages arXiv–2104.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2021. [Transformers as soft reasoners over language](#). In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3882–3890.
- Leyang Cui, Yu Wu, Jian Liu, Sen Yang, and Yue Zhang. 2021. [Template-based named entity recognition using BART](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1835–1845, Online. Association for Computational Linguistics.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Avia Efrat and Omer Levy. 2020. [The turking test: Can language models understand instructions?](#) *arXiv preprint arXiv:2010.11982*.
- Himanshu Gupta, Neeraj Varshney, Swaroop Mishra, Kuntal Kumar Pal, Saurabh Arjun Sawant, Kevin Scaria, Siddharth Goyal, and Chitta Baral. 2023. [“john is 50 years old, can his son be 65?” evaluating NLP models’ understanding of feasibility](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 407–417, Dubrovnik, Croatia. Association for Computational Linguistics.
- Tanmay Gupta, Amita Kamath, Aniruddha Kembhavi, and Derek Hoiem. 2021. [Towards general purpose vision systems](#). *arXiv preprint arXiv:2104.00743*.
- Christopher Hahn, Frederik Schmitt, Jens U. Kreber, Markus Norman Rabe, and Bernd Finkbeiner. 2021. [Teaching temporal logics to neural networks](#). In *International Conference on Learning Representations*.
- Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, et al. 2022. [FOLIO: Natural language reasoning with first-order logic](#). *arXiv preprint arXiv:2209.00840*.
- Peter Hase and Mohit Bansal. 2022. [When can models learn from explanations? a formal framework for understanding the roles of explanation data](#). In *Proceedings of the First Workshop on Learning with Natural Language Supervision*, pages 29–39, Dublin, Ireland. Association for Computational Linguistics.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. [Mistral 7b](#). *arXiv preprint arXiv:2310.06825*.
- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. 2020. [How can we know what language models know?](#) *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Pratik Joshi, Somak Aditya, Aalok Sathe, and Monojit Choudhury. 2020. [TaxiNLI: Taking a ride up the NLU hill](#). In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 41–55, Online. Association for Computational Linguistics.

- Daniel Khashabi. 2019. *Reasoning-Driven Question-Answering for Natural Language Understanding*. University of Pennsylvania.
- Kirby Kuznia, Swaroop Mishra, Mihir Parmar, and Chitta Baral. 2022. [Less is more: Summary of long instructions is better for program synthesis](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4532–4552, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Teven Le Scao and Alexander M Rush. 2021. How many data points is a prompt worth? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2627–2636.
- Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. [The Winograd Schema Challenge](#). In *Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning, KR’12*, pages 552–561. AAAI Press, Rome, Italy.
- Vladimir Lifschitz. 1989. Benchmark problems for formal nonmonotonic reasoning: Version 2.00. In *Non-Monotonic Reasoning: 2nd International Workshop Grassau, FRG, June 13–15, 1988 Proceedings 2*, pages 202–219. Springer.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shrutu Bhosale, Jingfei Du, et al. 2021. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021a. [Logiqa: a challenge dataset for machine reading comprehension with logical reasoning](#). In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3622–3628.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *arXiv preprint arXiv:2107.13586*.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. [Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.
- Man Luo, Sharad Saxena, Swaroop Mishra, Mihir Parmar, and Chitta Baral. 2022. [Biotabqa: Instruction learning for biomedical table question answering](#). *arXiv preprint arXiv:2207.02419*.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. [The SICK \(Sentences Involving Compositional Knowledge\) dataset for relatedness and entailment](#).
- R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference](#).
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2021. [Metaicl: Learning to learn in context](#). *arXiv preprint arXiv:2110.15943*.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2022a. [Reframing instructional prompts to GPTk’s language](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 589–612, Dublin, Ireland. Association for Computational Linguistics.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. [Cross-task generalization via natural language crowdsourcing instructions](#). *arXiv preprint arXiv:2104.08773*.
- Swaroop Mishra, Arindam Mitra, Neeraj Varshney, Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and Ashwin Kalyan. 2022b. [NumGLUE: A suite of fundamental yet challenging mathematical reasoning tasks](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3505–3523, Dublin, Ireland. Association for Computational Linguistics.
- Swaroop Mishra and Elnaz Nouri. 2023. [HELP ME THINK: A simple prompting strategy for non-experts to create customized content with models](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11834–11890, Toronto, Canada. Association for Computational Linguistics.
- Mutsumi Nakamura, Santosh Mashetty, Mihir Parmar, Neeraj Varshney, and Chitta Baral. 2023. [LogicAttack: Adversarial attacks for evaluating logical consistency of natural language inference](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13322–13334, Singapore. Association for Computational Linguistics.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. [Training language models to follow instructions with human feedback](#). *Advances in neural information processing systems*, 35:27730–27744.
- Mihir Parmar, Swaroop Mishra, Mirali Purohit, Man Luo, Murad Mohammad, and Chitta Baral. 2022. [In-BoXBART: Get instructions into biomedical multi-task learning](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 112–128, Seattle, United States. Association for Computational Linguistics.



- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. [Are NLP models really able to solve simple math word problems?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
- Pruthvi Patel, Swaroop Mishra, Mihir Parmar, and Chitta Baral. 2022. [Is a question decomposition unit all we need?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4553–4569, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mohammed Saeed, Naser Ahmadi, Preslav Nakov, and Paolo Papotti. 2021. [RuleBERT: Teaching soft rules to pre-trained language models.](#) In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1460–1476, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *ICLR*.
- Abulhair Saparov and He He. 2023. [Language models are greedy reasoners: A systematic formal analysis of chain-of-thought.](#) In *The Eleventh International Conference on Learning Representations*.
- Timo Schick and Hinrich Schütze. 2021. [Exploiting cloze-questions for few-shot text classification and natural language inference.](#) In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. 2019. [CLUTRR: A diagnostic benchmark for inductive reasoning from text.](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4506–4515, Hong Kong, China. Association for Computational Linguistics.
- Oyvind Tafjord, Peter Clark, Matt Gardner, Wen-tau Yih, and Ashish Sabharwal. 2019a. [Quarel: A dataset and models for answering questions about qualitative relationships.](#) In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7063–7071.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. [ProofWriter: Generating implications, proofs, and abductive statements over natural language.](#) In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634, Online. Association for Computational Linguistics.
- Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019b. [QuaRTz: An open-domain dataset of qualitative relationship questions.](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5941–5946, Hong Kong, China. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge.](#) In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jidong Tian, Yitian Li, Wenqing Chen, Liqiang Xiao, Hao He, and Yaohui Jin. 2021. [Diagnosing the first-order logical reasoning ability through LogicNLI.](#) In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3738–3747, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrusti Bhosale, et al. 2023. [Llama 2: Open foundation and fine-tuned chat models.](#) *arXiv preprint arXiv:2307.09288*.
- Bertie Vidgen, Dong Nguyen, Helen Margetts, Patricia Rossini, and Rebekah Tromble. 2021. [Cad: the contextual abuse dataset.](#)
- Liwen Wang, Rumei Li, Yang Yan, Yuanmeng Yan, Sirui Wang, Wei Wu, and Weiran Xu. 2022a. [Instructioner: A multi-task instruction-based generative framework for few-shot ner.](#) *arXiv preprint arXiv:2203.03903*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022b. [Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP](#)



- tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *ICLR*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Orion Weller, Nicholas Lourie, Matt Gardner, and Matthew E. Peters. 2020. [Learning from task descriptions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1361–1375, Online. Association for Computational Linguistics.
- Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–22.
- Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and Johan Bos. 2019. [Can neural networks understand monotonicity reasoning?](#) In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 31–40, Florence, Italy. Association for Computational Linguistics.
- Qinyuan Ye and Xiang Ren. 2021. [Learning to generate task-specific adapters from task description](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 646–653, Online. Association for Computational Linguistics.
- Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. [Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3914–3923, Hong Kong, China. Association for Computational Linguistics.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.
- Weihaoyu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. [Reclor: A reading comprehension dataset requiring logical reasoning](#). In *International Conference on Learning Representations*.
- Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang, and Guy Van Den Broeck. 2023. On the paradox of learning to reason from data. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 3365–3373.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein. 2021. [Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2856–2878, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. [“going on a vacation” takes longer than “going for a walk”](#): A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3363–3369, Hong Kong, China. Association for Computational Linguistics.

## A Example Prompt for Sentence Generation

Below is the general prompt structure prompted to generate data. The prompt schema, as depicted in Figure 3, comprise three crucial components:

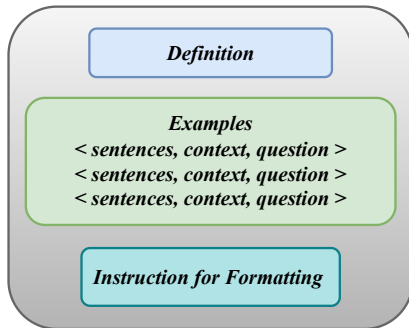


Figure 3: Schematic representation of prompt.

**Definition** provides a detailed explanation of the task and offers a natural language representation of the reasoning pattern for which we are generating sentences.

**Examples** provide sample sentences that need to be generated. We also illustrate how these sentences will be utilized in later stages, emphasizing the importance of coherence and the inclusion of relevant ontological concepts.

**Format** We provide specific formatting instructions to guide the generation of sentences.

Figure 4 illustrates an example prompt for the inference rule, namely, ‘modus tollens’ from propositional logic (PL). Modus tollens is formally represented as  $((p \rightarrow q) \wedge \neg q) \vdash \neg p$ , which can be understood in natural language as “If  $p$  implies  $q$ , and we know  $\neg q$ , then we can conclude  $\neg p$ .” In this prompt, the definition provides a comprehensive description of the inference rule in natural language. To encourage the generation of more relevant and coherent sentences, the prompt includes an examples section that demonstrates how the generated sentences will be utilized in a later stage. This serves, as an illustration, to guide GPT-3 in producing suitable outputs. In Figure 4, we present three examples involving sentences  $p$  and  $q$ , along with their respective contexts and questions. The prompt also includes instructions on how the generated sentences should be formatted.

## B Extended Related Work

As LLMs such as GPT-4, and Bard continue to evolve rapidly, it becomes increasingly crucial to

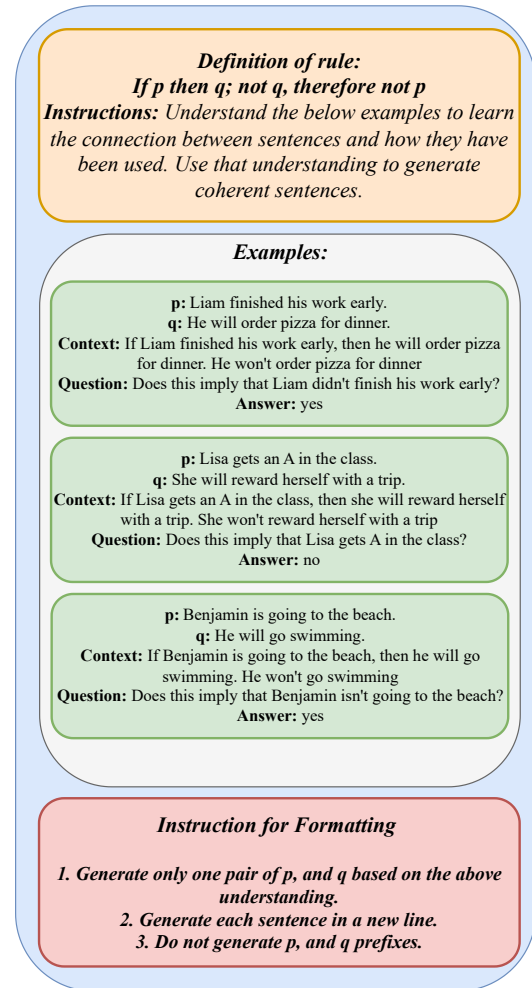


Figure 4: Example prompt for *Modus Tollens* inference rule from PL.

evaluate their diverse language capabilities, as well as those of forthcoming LLMs. Recently, many datasets have been created that evaluate different language understanding skills such as pronoun resolution (Sakaguchi et al., 2021; Levesque et al., 2012), commonsense reasoning (Talmor et al., 2019), numerical reasoning (Dua et al., 2019; Patel et al., 2021; Mishra et al., 2022b), qualitative reasoning (Tafjord et al., 2019b,a), temporal reasoning (Zhou et al., 2019), and feasibility reasoning (Gupta et al., 2023). Now, we present the advancements in prompt and instruction tuning using LLMs.

**Prompt Learning** The introduction of LLMs has significantly shifted the research trend in NLP to prompt-based learning methodologies (Liu et al., 2021b). Many studies have been conducted to investigate the efficacy of prompt-based learning in various applications including Text classification (Yin et al., 2019), Natural Language Inference

(NLI) (Schick and Schütze, 2021), and Question Answering (QA) (Jiang et al., 2020), Information Extraction (IE) (Chen et al., 2021; Cui et al., 2021), to name a few. In a recent development, the T0 model employs prompts to achieve zero-shot generalization across various NLP tasks (Sanh et al., 2021). Le Scao and Rush (2021) suggested that the use of prompts could be as valuable as hundreds of data points on average.

**Instruction Learning** Efrat and Levy (2020) was focused on whether existing LLMs understand instructions. The same work in the field of instruction by Hase and Bansal (2022); Ye and Ren (2021); Gupta et al. (2021); Zhong et al. (2021) has been proposed to show that models follow natural language instructions. In addition, Weller et al. (2020) developed a framework focusing on NLP systems that solve challenging new tasks based on their description. Mishra et al. (2021) have proposed natural language instructions for cross-task generalization of LLMs. Similarly, PromptSource (Sanh et al., 2021) and FLAN (Wei et al., 2021) were built for leveraging instructions and achieving zero-shot generalization on unseen tasks. Moreover, Parmar et al. (2022) shows the effectiveness of instructions in multi-task settings for the biomedical/clinical domain. Furthermore, Mishra et al. (2022a) discussed the impact of task instruction reframing. Min et al. (2021) introduced a framework to better understand in-context learning. Ouyang et al. (2022) proposed the InstructGPT model that is fine-tuned with human feedback to follow instructions. Wang et al. (2022a) has developed an instruction-based multi-task framework for few-shot Named Entity Recognition (NER) tasks. In addition, many approaches have been proposed to improve model performance using instructions (Wu et al., 2022; Lin et al., 2021; Wang et al., 2022b; Luo et al., 2022; Kuznia et al., 2022; Patel et al., 2022; Mishra and Nouri, 2023).

**Logic and NLI Datasets** FraCas (Bernardy and Chatzikyriakidis, 2020) offers a unique approach to temporal semantics by converting syntax trees into logical formulas tailored for inference, emphasizing temporal elements such as references, adverbs, aspectual classes, and progressives. The Monotonicity Entailment Dataset (MED) (Yanaka et al., 2019) dives deep into monotonicity reasoning within NLI, probing the synergy between lexical and syntactic structures and spotlighting inherent challenges in both upward and downward monotonic reasoning trajectories. The SICK

(Marelli et al., 2014) dataset, with its foundation in 10,000 English sentence pairs, is designed to rigorously evaluate semantic relatedness and entailment, leveraging crowdsourced annotations for precision. HANS, or Heuristic Analysis for NLI Systems (McCoy et al., 2019), stands out by rigorously scrutinizing the dependability of NLI models, putting the spotlight on potential pitfalls tied to syntactic heuristics such as lexical overlap. Lastly, CAD (Vidgen et al., 2021) introduces a meticulously crafted dataset from Reddit entries, targeting the detection of online abuse. Nakamura et al. (2023) introduced LogicAttack, a method for performing adversarial attacks on NLI models using PL inference rules and proposed dataset of  $\sim 9k$  attack samples derived from the SNLI dataset. In contrast to these works, *LogicBench* evaluates the logical reasoning capabilities of LLMs beyond the NLI task, focusing on the application of individual inference rules.

## C Examples of Data Instances

This section provides examples of (*context, question, answer*) tuples corresponding to each inference rule and reasoning pattern. Additionally, it highlights the diverse range of question variations within the dataset associated with each inference rule and reasoning pattern.

### C.1 Word Cloud

Figure 5 provides a word cloud derived from the *LogicBench(Eval)*. This word cloud highlights the logical nature and diversity of our evaluation dataset. Words such as ‘if’, ‘normally’, ‘usually’, and ‘then’ are prominently featured, suggesting their frequent use in the dataset, and suggesting the logical nature of the dataset. Moreover, we can also observe several words consisting of different ontologies such as ‘cat’, ‘car’, ‘garden’, and many more, suggesting diversity in the dataset.

### C.2 Propositional Logic (PL)

Here, we discuss examples of each inference rule present in the PL of the *LogicBench* as shown in Table 9. Table 9 has context related to the inference rule and different variations of the question according to the rule. For instance, the first row of Table 9 shows the example for inference rule, *Hypothetical Syllogism (HS)*, formally expressed as  $((p \rightarrow q) \wedge (q \rightarrow r)) \vdash (p \rightarrow r)$ . The context represents the premise, i.e.,  $((p \rightarrow q) \wedge (q \rightarrow r))$ ,





Rule	Context	Question
HS	If Jim cleans his room, he will receive a reward and if he receives a reward, he will use it to buy a new toy. So, Jim decided to tidy up his room, hoping to earn a reward. He diligently gathered his clothes, organized his toys, and dusted every surface. After a few hours of hard work, Jim's room was spotless.	<b>Q1:</b> If Jim cleaned his room, does this imply that he will buy a new toy? (Yes) <b>Q2:</b> If Jim didn't clean his room, does this entail that he won't buy a new toy? (No) <b>Q3:</b> If Jim cleaned his room, does this imply that he won't buy a new toy? (No) <b>Q4:</b> If Jim didn't clean his room, does this imply that he will buy a new toy? (No)
DS	Either Chloe is studying for her exams, or Mila is going on vacation, or both scenarios are unfolding. It is unclear which of the two options is true at this point. However, one thing is certain - Chloe is not studying for her exams.	<b>Q1:</b> Does this entail that Mila is going on vacation? (Yes) <b>Q2:</b> Does this mean that Mila isn't going on vacation? (No)
CD	If I decide to go for a walk, I will be able to breathe in some fresh air and revitalize myself. On the other hand, if I choose to stay home, I will have the opportunity to enjoy a movie. One thing is certain, either I go for a walk or I stay home. It remains uncertain which of the two options I will ultimately choose. It is entirely possible that I might opt for the walk, or perhaps I will find myself drawn to staying home, or even both possibilities might come to fruition.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) I will get some fresh air or (b) I will watch a movie (Yes) <b>Q2:</b> (a) I won't get some fresh air and (b) I will watch a movie (No) <b>Q3:</b> (a) I will get some fresh air and (b) I won't watch a movie (No) <b>Q4:</b> (a) I won't get some fresh air and (b) I won't watch a movie (No)
DD	If I decide to order takeout, it means I will save time. On the other hand, if I choose to cook a meal, it means I will save money. The interesting thing is that I am in a situation where I won't be able to save time or money. It is uncertain whether I won't save time or I won't save money, or it could even be both scenarios. The only thing that is clear is that at least one of these possibilities is true.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) I don't order takeout or (b) I don't cook a meal (Yes) <b>Q2:</b> (a) I order takeout and (b) I cook a meal (No) <b>Q3:</b> (a) I don't order takeout and (b) I cook a meal (No) <b>Q4:</b> (a) I order takeout and (b) I don't cook a meal (No)
BD	If it is sunny outside, then I will go for a walk. However, if it rains, I will stay inside. Currently, it is uncertain whether it is raining or not, but I do know that at least one of the following is true: either it is raining or I will not go for a walk. It is possible that only one of these statements is true, or perhaps both are true.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) we will stay inside or (b) it is not sunny (Yes) <b>Q2:</b> (a) we will not stay inside and (b) it is sunny (No) <b>Q3:</b> (a) we will stay inside and (b) it is sunny (No) <b>Q4:</b> (a) we will not stay inside and (b) it is not sunny (No)
MT	If Mason decides to leave his job, he will not receive any salary. However, against all odds, Mason still receives his salary. He finds himself receiving his regular paycheck.	<b>Q1:</b> Does this infer that Mason didn't leave his job? (Yes) <b>Q2:</b> Does this infer that Mason left his job? (No)
MI	Rohan woke up in the morning and realized that he had forgotten his lunch. Knowing that if he forgets his lunch, he will not eat at school, he felt disappointed.	Based on context, can we say, at least one of the following must always be true? <b>Q1:</b> (a) Rohan didn't forget his lunch and (b) he will not eat at school (Yes) <b>Q2:</b> (a) Rohan forgot his lunch and (b) he will eat at school (No) <b>Q3:</b> (a) Rohan forgot his lunch and (b) he will not eat at school (No) <b>Q4:</b> (a) Rohan didn't forget his lunch and (b) he will eat at school (No)
CT	At least one of two things is true about Tom - he is either an avid reader or he devours books of all genres. We are unsure which one of these statements is true or if both are true. It could be that only the first statement is true, or only the second statement is true, or even that both are true.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) he devours books of all genres or (b) Tom is an avid reader (Yes) <b>Q2:</b> (a) he doesn't devour books of all genres and (b) Tom is an avid reader (No) <b>Q3:</b> (a) he devours books of all genres and (b) Tom isn't an avid reader (No) <b>Q4:</b> (a) he doesn't devour books of all genres and (b) Tom isn't an avid reader (No)

Table 9: Examples of context and question-answer pairs for each rule of *Proportional logic* from the LogicBench; HS: Hypothetical Syllogism, DS: Disjunctive Syllogism, CD: Constructive Dilemma, DD: Destructive Dilemma, BD: Bidirectional Dilemma, MT: Modus Tollens, MI: Material Implication, CT: Commutation

Rule	Context	Question
UI	All students are required to take an examination in order to fulfill the requirements for their degree. Reema, being a student, is also expected to fulfill the requirements.	<b>Q1:</b> Does Reema need to take an exam to complete her degree? (Yes) <b>Q2:</b> Does Reema need not to take an exam to complete her degree? (No)
EG	The marathon race was won by James, who emerged as the champion.	<b>Q1:</b> Does this imply that someone won the marathon race? (Yes) <b>Q2:</b> Does this mean that no one won the marathon race? (No)
MP	If someone is extremely tired, then they will seek some rest and relaxation. Today, Jack finds himself utterly exhausted.	<b>Q1:</b> Does this entail that he will take rest? (Yes) <b>Q2:</b> Does this entail that he won't take rest? (No)
HS	If all the necessary supplies have been purchased by someone, then they can initiate the project. Once the project is started by someone, they will complete it within the expected time-frame.	<b>Q1:</b> If Lily bought all the necessary supplies, does this mean that she will finish it on time? (Yes) <b>Q2:</b> If Lily didn't buy all the necessary supplies, does this imply that she won't finish it on time? (No) <b>Q3:</b> If Lily bought all the necessary supplies, does this entail that she won't finish it on time? (No) <b>Q4:</b> If Lily didn't buy all the necessary supplies, does this imply that she will finish it on time? (No)
DS	It is known that one of the following options is true: someone goes to a museum or someone visits a park. The specific scenario could involve only the option to go to a museum being true, or only the option to visit a park being true, or both options being true. However, it is stated that Jill is unable to go to a museum.	<b>Q1:</b> Does this imply that she can visit a park? (Yes) <b>Q2:</b> Does this entail that she can't visit a park? (No)
CD	If someone is painting a picture, then they will frame it. Similarly, the natural course of action for a writer would be to publish their completed story. In this scenario, it is certain that at least one of the following statements holds true: (1) John is currently engrossed in painting a picture, or (2) John is currently immersed in the act of writing a story. It should be emphasized that we are unaware of which statement specifically applies, as there is a possibility that either (1) alone is true, or (2) alone is true, or even that both (1) and (2) are simultaneously true.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) he will frame it and (b) he will publish it. (Yes) <b>Q2:</b> (a) he won't frame it and (b) he will publish it. (No) <b>Q3:</b> (a) he will frame it and (b) he won't publish it. (No) <b>Q4:</b> (a) he won't frame it and (b) he won't publish it. (No)
DD	If someone is taking care of their health, then they will be fit. However, indulging in unhealthy habits can make individuals susceptible to various diseases. The truth is, we can be certain about at least one of the following possibilities: either Jenny won't be fit and healthy, or she won't be prone to diseases. It is important to note that we are unaware of which statement is accurate. It could be the case that only the first statement is true, only the second statement is true, or both statements hold validity.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) Jenny doesn't take care of her health and (b) she doesn't indulge in unhealthy habits (Yes) <b>Q2:</b> (a) Jenny takes care of her health and (b) she indulges in unhealthy habits (No) <b>Q3:</b> (a) Jenny doesn't take care of her health and (b) she indulges in unhealthy habits (No) <b>Q4:</b> (a) Jenny takes care of her health and (b) she doesn't indulge in unhealthy habits (No)
BD	If an individual consumes a significant amount of water, they will experience a state of hydration. Conversely, if excessive amounts of sugar are ingested, a sugar crash will ensue. It is known that at least one of the following statements is true: either the Jane consumes ample water or she will not experience a sugar crash. However, the actual veracity of either statement remains ambiguous, as it could be the case that only the first statement is true, only the second statement is true, or both statements are true.	Can we say at least one of the following must always be true? <b>Q1:</b> (a) she will feel hydrated and (b) she doesn't eat too much sugar (Yes) <b>Q2:</b> (a) she won't feel hydrated and (b) she eats too much sugar (No) <b>Q3:</b> (a) she will feel hydrated and (b) she eats too much sugar (No) <b>Q4:</b> (a) she won't feel hydrated and (b) she doesn't eat too much sugar (No)
MT	If someone decides to go to the park, it is required that they wear a mask. However, in this particular situation, John does not wear a face covering.	<b>Q1:</b> Does this imply that John doesn't visit the park? (Yes) <b>Q2:</b> Does this entail that John visits the park? (No)

Table 10: Examples of context and question-answer pairs for each rule of *First order logic* from the LogicBench.

Rule	Context	Question
DRI	Once upon a time, in a land filled with animals, there were two popular mammalian creatures, cats and dogs. Mammals typically possessed a coat of fur, which kept them warm and protected. However, cats were an exception to this rule, as their bodies lacked fur. Nonetheless, both cats and dogs were beloved by many for their unique traits. Dogs, known for their loyalty, were particularly cherished by humans.	<b>Q1:</b> Does this imply that dogs have fur? (Yes) <b>Q2:</b> Does this entail that dogs don't have fur? (No)
DRS	John and Mary were expecting their first child, filled with the anticipation and excitement that all parents feel. Parents are usually loving and supportive. Parents are normally responsible. However, something seemed amiss in their relationship. Mary, usually affectionate and caring, seemed distant and uninvolved. On the other hand, John, known for his responsible nature, started neglecting his duties and became unreliable.	<b>Q1:</b> Does this imply that Mary is responsible and John is loving and supportive? (Yes) <b>Q2:</b> Does this entail that Mary isn't responsible and John is loving and supportive? (No) <b>Q3:</b> Does this imply that Mary is responsible and John isn't loving and supportive? (No) <b>Q4:</b> Does this entail that Mary isn't responsible and John isn't loving and supportive? (No)
DRD	Jenny and Anna are known for their tall stature, which is often associated with playing basketball. However, Anna might be an exception to this norm.	<b>Q1:</b> Does this entail that Jenny plays basketball? (Yes) <b>Q2:</b> Does this mean that Jenny doesn't play basketball? (No)
DRO	In the bird kingdom, there are many different species that possess unique characteristics. One such species is the hummingbird, known for its ability to hover in mid-air and its vibrant colors. While most birds engage in the annual migration south for the winter, the hummingbird chooses to stay put and brave the cold weather. This decision sets the hummingbird apart from its fellow avian companions, as it relies on its resilience and resourcefulness to survive the harsh conditions.	<b>Q1:</b> Does this mean that all other birds than hummingbirds migrate south for the winter? (Yes) <b>Q2:</b> Does this mean that all other birds than hummingbirds don't migrate south for the winter? (No)
RE1	In a world where animals are often regarded as intelligent creatures, there is a captivating tale that revolves around cats, dogs, and horses. It is commonly believed that most animals possess a level of intellect. However, there is an intriguing twist to this belief as it is known that either cats or dogs are not considered particularly intelligent. As the story unfolds, we delve into the lives of these remarkable creatures, their interactions, and the unique qualities that each of them possesses.	<b>Q1:</b> Does this entail that horses are considered to be intelligent creatures and exactly one of the cats or dogs is not considered intelligent? (Yes) <b>Q2:</b> Does this mean that horses aren't considered to be intelligent creatures and exactly one of cats or dogs is not considered intelligent? (No) <b>Q3:</b> Does this mean that horses are considered to be intelligent creatures and exactly one of cats or dogs is considered intelligent? (No) <b>Q4:</b> Does this implies that horses aren't considered to be intelligent creatures and exactly one of cats or dogs is considered intelligent? (No)
RE2	In the realm of cat communication, meowing serves as a fundamental aspect of their vocal repertoire. However, intriguingly enough, there exists a distinct species of cat that deviates from this conventional norm. This peculiar feline defies the expectations associated with its kind by refraining from emitting any meows whatsoever.	<b>Q1:</b> Does this entail that exactly one species of cat doesn't meow? (Yes) <b>Q2:</b> Does this imply that exactly one species of cat meows? (No)
RE3	In a world where cars were known for having four wheels, it was considered a common fact that wheels typically came equipped with spokes. However, amidst this widespread understanding, there was an exception. At least one wheel defied this norm and stood out from the rest by not having any spokes at all.	<b>Q1:</b> Does this imply that cars have four wheels with spokes? (Yes) <b>Q2:</b> Does this mean that cars don't have four wheels with spokes? (No)
RAP	In the midst of a heated argument, John adamantly claims that Sally was present at the store. However, Jane strongly opposes John's assertion, insisting that Sally was indeed absent from the store.	<b>Q1:</b> If John's evidence is more reliable than Jane's, does this mean that Sally was in the store? (Yes) <b>Q2:</b> If John's evidence is more reliable than Jane's, does this mean that Sally wasn't in the store? (No) <b>Q3:</b> If John's evidence is less reliable than Jane's, does this entail that Sally was in the store? (No) <b>Q4:</b> If John's evidence is less reliable than Jane's, does this imply that Sally wasn't in the store? (Yes)

Table 11: Examples of context and question-answer pairs for each rule of *Non-monotonic logic* from the LogicBench.

Each prompt is shown below for PL, FOL, and NM templated context to narrative conversion.

### 1. Prompt for templated to narrative conversion for PL(MT):

**Rule:** Condition: If p then q; **Situation:** not q

**Context:** Templated context

Improve the context in human-like language and make story with rephrased sentences.

**Instructions for generating a good story:**

1. While generating a story, Use the rephrased sentences from the context of the story.
2. Make sure to include sentences corresponding to the condition and situation from the rule in the story.
3. Do not add any other extra information.
4. To generate a story, DO NOT change the name of the main character from the context, if there's any. If the main character within the context is 'I', please generate a story with 'I' as the main character.
5. Generate only one paragraph with rephrased sentences.

### 2. Prompt for templated to narrative conversion for FOL:

**Context:** Templated context

Improve the context in human-like language and make a rephrase sentences around it.

**Instructions for generating a good story:**

1. Make sure to include only rephrased sentences in the story.
2. Do not add any other extra information and do not reach any conclusion based on context.
3. To generate a story, use only generic sentences without mentioning any main character.
4. Extract a specific case sentence from the context, rephrase it, and append it at the end of the context only.
5. Generate only one paragraph with rephrased sentences.
6. DO NOT add any prefixes.

### 3. Prompt for templated to narrative conversion for NM:

**Context:** Templated context

Improve the context in human-like language and rephrase sentences around it.

**Instructions for generating a good story:**

1. Make sure to include only rephrased sentences in the story.
2. Do not add any other extra information, but add some story-like context.
3. To generate a story, DO NOT change the name of the main character from the context, if there's any. Use the same names/words for main character. If the main character within the context is 'I', please generate a story with 'I' as the main character.
4. Generate only one paragraph with rephrased sentences.

To ensure the quality of the narrative in the *LogicBench(Eval)* for task 1 and task 2, we have created category-specific prompts to convert the templated context to convert more human-like narrative. In total, we have created three different prompts each for PL[1D], FOL[2D], and NM[3D]. The prompts are designed to ensure that the logical connection is established in the narrative. Each rule of PL, FOL, and NM has a unique logical progression that should be followed in the narrative is the motivation for us to go with three different instruction-based prompts. In the prompt, used for PL rules, the logical rule is also mentioned in two parts "Condition" and "Situation". Consider the example of logical rule MT - "If p then q; not q; therefore not p", the templated context will have "If p then q; not q;" so condition will be the first part "If p then q" while the situation will be the second part "not q". In other terms, the specific conditions are the rules to be followed while the situation is the case present in the logical rule. In contrast, for FOL, we do not need to focus more on the specific condition rather we have to make sure a generalized case is present in the narrative with 1-2 specific sentences related to the rule. For NM, we do not have such a rule instead we have a logical connection between sentences and hence we go with the instruction-based prompting.

## E Task Instance Generation

This section discusses the task instance generation step, which is the last step in Fig. 2 in detail. The overall section is comprised of two subsections for BQA and MCQA. The subsection related to BQA



discusses the variation generation step in detail with question generation while the MCQA subsection provides details of prompts used for incorrect option generation along with five pre-defined sets of questions.

### E.1 Task 1- BQA (variation generation)

As seen in Fig. 2, BQA has a narrative along with question, and answer pairs with variation in question and corresponding answer. As discussed in the main paper Data creation 3.2, the narrative is created using the pre-defined rules for PL, FOL, and NM while a question is asked based on what can be entailed from the given context. For example, HS can be defined as "If p then q; if q then r; therefore, if p then r" and the narrative will have "If p then q; if q then r" and the question is asked as "Can we conclude if p then r?" having answer as "yes". Now, other variations in the question are asked in the following ways by negating the sentence (p and r) comb: Variation 1: "Can we conclude if  $\neg$ p then r?", Variation 2: "Can we conclude if p then  $\neg$ r?", Variation 3: "Can we conclude if  $\neg$ p then  $\neg$ r?". If there's only one axiom (p) present in the question then there are only a 2 variations that can be made asking about (p) and ( $\neg$ p).

### E.2 Task 2-MCQA (Question selection and incorrect option generation)

We have MCQA as task 2 in LogicBench. In this, we have one correct option from the four options, and three options are incorrectly generated using prompting. The question is a bit different from the BQA question formation as we have MCQA and have to identify which conclusion can be derived from the narrative. The question is randomly selected from the pre-defined set of questions. Here, the five different questions are as follows:

1. What would be the most appropriate conclusion based on the given context?
2. Considering the provided context, what conclusion would be deemed most suitable?
3. In light of the context provided, what conclusion can be considered the most appropriate?
4. Based on the context, what conclusion would be deemed most suitable?
5. Taking into account the context provided, what conclusion would be most appropriate?

The narrative for MCQA is the same as the BQA's narrative and there's no change. Comparing BQA methodology with MCQA, MCQA's correct option is the question asked in the BQA which can be concluded from the logical rule present in the narrative. For example, HS can be defined as "If p then q; if q then r; therefore, if p then r" and the narrative will have "If p then q; if q then r" and the question is randomly chosen from the set of questions as mentioned above and the correct option from the different multiple choice will be "if p then r". Based on the information present, we have generated an incorrect option using prompting [4E.2]. Refer this link for more examples - <https://anonymous.4open.science/r/LogicBench-EEBB>.

### 4. Prompt for incorrect option generation:

**Context:** Narrative

**Question:** Randomly chosen question

**Correct option:** Conclusion from narrative

Create three more irrelevant, incorrect options to create multiple choice. Sometime use different persons name as well in creating incorrect options.

**Instruction to generate options:**

1. Generate only three unique choices each in new line.
2. Do not use any prefix while generating.

## F Mitigation of Errors in LogicBench

While validating *LogicBench(Eval)*, we encountered errors within the synthetically generated narrative-based context. We mitigate these errors manually, categorizing them into two groups: (i) eliminating leakage of logical conclusions and (ii) ensuring the inclusion of intended logical premises. In the first category, we found  $\sim 15\%$  of the narrative-based contexts (out of 500 total instances) were found to explicitly present the logical conclusion as a response to the question, bypassing the logical derivation process. This enables the model to extract the final logical conclusion from the context rather than derive it logically. In the second category, we found  $\sim 8\%$  of the narrative-based contexts (out of 500 total instances) where narration lacked some necessary premise sentences crucial for reaching a logical conclusion. To address this issue, we manually incorporated those sentences to ensure the quality of the data. For the MCQA task, during the generation of three incor-

rect options, we found instances where the model produces two semantically similar options, resulting in the creation of ambiguous choices in  $\sim 11\%$  of the cases (out of a total of 500 instances). We manually mitigated all these errors from the data instances ensuring the high quality of our validation data. We believe that these two versions aim to accommodate different evaluations to explore the logical reasoning capabilities of LLMs.

## G Few-shot Experiments

This section discusses the performance of the different LLMs in a few-shot setting on the *LogicBench(Eval)<sub>BQA</sub>*. Here, we only present a case study on the BQA task. For the fair comparison with Table 6, we analyze an average performance across  $A(Yes)$ . Table 15 shows the performance for each inference rule and reasoning patterns achieved by Llama-2, Mistral, Gemini, ChatGPT, and GPT-4.

As suggested in (Lu et al., 2022), prompting models are sensitive to in-context examples. Hence, we see mixed performance in Table 15 across all models. From Table 15, we can observe that in-context examples are helpful for Llama-2 since it consistently outperforms zero-shot baselines by large margins in terms of  $A(Yes)$ . Llama-2 is remarkably good at following the in-context exemplars and mimicking the process to reach the correct conclusions. Thus, leveraging the in-context exemplars, Llama-2 achieves high accuracy in a few-shot setting. Specifically, Mistral consistently shows degraded performance for all logic types. However, ChatGPT improves performance on NM reasoning. Improved performance in NM reasoning demonstrates that the inclusion of in-context examples enhances the ability of these models to comprehend the nuanced meanings of logical terms such as “usually” and “typically”. In particular, we see that Gemini and GPT-4 improve performance on PL and FOL, respectively, but show competitive performance on NM.

## H Experimental Setup

### H.1 Extended Discussion on Experiments

**Zero-shot setting** We evaluate GPT-4, and ChatGPT (GPT-3.5-Turbo) by utilizing their APIs provided by OpenAI<sup>3</sup>. We evaluate Google Gemini-Pro by utilizing its API provided by Google<sup>4</sup> The

<sup>3</sup><https://platform.openai.com/docs/guides/gpt>

<sup>4</sup><https://ai.google.dev/>

evaluation is conducted on the versions of GPT-4, ChatGPT, and Gemini released in January 2024. It’s important to note that these models are regularly updated, so when reproducing the results presented in Table 6 and Table 7 (main paper), there is a possibility of variations. For Llama-2 and Mistral, we utilize the 7B-Chat, and 7B-Instruct-v0.2 versions, respectively, from the huggingface model repository<sup>5</sup>.

**Experiments on other logic datasets** In single and multi-task experiments on other logic datasets, we fine-tune the T5-large model for 10 epochs with a batch size of 16, 1024 maximum input length, an adaptive learning rate of  $5e - 05$ , and an AdamW optimizer for each experiment. All experiments are performed using NVIDIA RTX A6000 GPUs.

### H.2 Prompts

All the experiments conducted in the zero-shot setting were performed using three distinct prompts. The reported results in Table 6 (main paper) represent the average performance across these prompts. All the prompts follow the common pattern which includes task description and formatting instructions. The following are the three different prompts utilized in the experiments:

#### Prompt 1:

Given the context and question, think step-by-step logically to answer the question. Answer the question ONLY in ‘yes’ or ‘no’. Please use the below format:  
**Context:** [text with logical rules]  
**Question:** [question based on context]  
**Reasoning steps:** [generate step-by-step reasoning]  
**Answer:** Yes/No

#### Prompt 2:

Let’s think step-by-step to answer the question given context. Answer the question ONLY in ‘yes’ or ‘no’. Please use the below format:  
**Context:** [text with logical rules]  
**Question:** [question based on context]  
**Reasoning steps:** [generate step-by-step reasoning]  
**Answer:** Yes/No

<sup>5</sup><https://huggingface.co/models>

### Prompt 3:

Given the context that contains rules of logical reasoning in natural language and question, perform step-by-step reasoning to answer the question. Based on context and reasoning steps answer the question ONLY in 'yes' or 'no'. Please use the below format:

**Context:** [text with logical rules]

**Question:** [question based on context]

**Reasoning steps:** [generate step-by-step reasoning]

**Answer:** Yes/No

For MCQA task, the reported results in Table 7 (main paper) represent the average performance across three prompts. All the prompts follow the common pattern which includes task description and formatting instructions. The following are the three different prompts utilized in the experiments:

### Prompt 1:

Given the context that contains rules of logical reasoning in natural language, question, and options, perform step-by-step reasoning to answer the question. Answer the question ONLY in 'choice\_1' or 'choice\_2' or 'choice\_3' or 'choice\_4'. DO NOT generate any text with it. Based on context and reasoning steps, choose the correct option.

Please use the below format:

**Context:** [text with logical rules]

**Question:** [question based on context]

**Options:** [four different choices]

**Reasoning steps:** [generate step-by-step reasoning]

**Answer:** choice\_1 / choice\_2 / choice\_3 / choice\_4

### Prompt 2:

Let's think step-by-step to choose the correct option for the question given context. Answer the question ONLY in 'choice\_1' or 'choice\_2' or 'choice\_3' or 'choice\_4'. DO NOT generate any text with it. Based on context and reasoning steps, choose the correct option.

Please use the below format:

**Context:** [text with logical rules]

**Question:** [question based on context]

**Options:** [four different choices]

**Reasoning steps:** [generate step-by-step reasoning]

**Answer:** choice\_1 / choice\_2 / choice\_3 / choice\_4

### Prompt 3:

Given the context, question, and options, think step-by-step logically to answer the question by choosing the correct option. Answer the question ONLY in 'choice\_1' or 'choice\_2' or 'choice\_3' or 'choice\_4'. DO NOT generate any text with it. Based on context and reasoning steps, choose the correct option.

Please use the below format:

**Context:** [text with logical rules]

**Question:** [question based on context]

**Options:** [four different choices]

**Reasoning steps:** [generate step-by-step reasoning]

**Answer:** choice\_1 / choice\_2 / choice\_3 / choice\_4

## I Further Discussion on Results

**Effect on other logic datasets** Our experiments were carried out in two settings: single-task (fine-tuning and evaluation on one dataset) and multi-task (fine-tuning on all four datasets combined, with separate evaluations for each dataset). App. H describes a detailed experimental setup. Table 12 represents the accuracy comparison between LogicT5 and baseline T5-large in both single-task and multi-task settings.

The results indicate that training LLMs on *LogicBench(Aug)* has a greater impact on logic datasets that primarily focus on logical reasoning, such as FOLIO and LogicNLI. Hence, we can observe that LogicT5 consistently outperforms the baseline for LogicT5 and FOLIO. However, LogiQA and ReClor encompass other forms of reasoning in addition to logical reasoning, hence, LogicT5 demonstrates competitive performance on them.

**LLMs sometimes overlook contextual information** In this case study, we manually examined all reasoning chains for PL and NM generated by GPT-4. The example below illustrates the observed

Methods	Models	LogiQA	FOLIO	LogicNLI	ReClor
Single-Task	T5-large	16.8	69.6	82.3	35.4
	LogicT5	<b>16.9</b>	<b>71.2</b>	<b>84.4</b>	<b>36.8</b>
Multi-Task	T5-large	<b>21.8</b>	83.8	68.2	<b>42.8</b>
	LogicT5	19.7	<b>85.6</b>	<b>69.8</b>	40.0

Table 12: Performance comparison between LogicT5 and baseline T5-large in terms of accuracy.

pattern. Specifically, in step 5 of the reasoning chain, the model deviates from the provided context and introduces unrelated information, resulting in an incorrect prediction. Although these patterns are not prevalent, we found that  $\sim 1\%$  of the total analyzed reasoning chains exhibited this pattern in PL and NM for BQA, and  $\sim 0.8\%$  for MCQA.

**Rule:** Modus Tollens (PL)

**Context:** James had always been forgetful, especially when it came to wearing a sweater. his mother had repeatedly told him, 'If he forgets to wear a sweater, then he will catch a cold.' throughout the day, he went about his activities, oblivious to the dropping temperature. hours later, he returned home, but James didn't catch a cold.

**Question:** Does this infer that James didn't forget to wear a sweater?

**Gold Answer:** Yes

**Reasoning steps:**

1. James is known to be forgetful about wearing a sweater.
2. His mother stated that if he forgets to wear a sweater, he will catch a cold.
3. The temperature dropped, which would increase the likelihood of catching a cold if he were not wearing a sweater.
4. Despite the dropping temperature, James did not catch a cold.
5. The absence of a cold does not necessarily mean that James remembered to wear a sweater, as there could be other reasons why he did not catch a cold (e.g., he may have a strong immune system, spent little time outdoors, or was in a warm environment).
6. Therefore, the fact that James did not catch a cold does not logically infer that he remembered to wear a sweater.

Answer: No

'if' and 'then' in them". From the results, we can observe that 13 of the 20 sentences in response to the prompt were FOL kind and only 7 were propositional kind. This shows that LLMs' comprehend simple FOL sentences, thus, showing high overall accuracy on simpler FOL compared to PL. However, we interestingly observed that, as the depth increases, PL-based examples carrying a specific case/story become easier to understand compared to FOL-based generalized context.

**Negations are hard to understand when embedded with logical rules.** Here, Table 16 presents reasoning chains corresponding to PL and FOL supporting the effect of negation on their predictions.

**Performance of BQA vs. MCQA** By analyzing the reasoning chains of CT and DD for PL, we observe that LLMs often employ the elimination method in MCQA. They prioritize options to select the most appropriate conclusion, thereby improving performance. Conversely, for NM, the models typically struggle to differentiate between two options and tend to hallucinate by assuming information, resulting in lower performance. We present reasoning chains corresponding to PL and NM supporting their behavior on both BQA and MCQA tasks at <https://github.com/Mihir3009/LogicBench>.

**Results on Yi-34B** Table 17 provides results for both BQA and MCQA tasks on *LogicBench*.

**PL vs. FOL** To get an indication of this we gave GPT-4 prompt - "Give twenty statements that have



Rule	Generate Sentences in Step 1	NL logical expressions
MT	p: Liam finished his work early. ~p: Liam did not finish his work early. q: He will order pizza for dinner. ~q: He will not order pizza for dinner.	<b>Context:</b> Liam had finished his work early for the day, which meant that he would typically have ordered pizza for dinner. However, on this particular day, he decided against ordering pizza and opted for something else instead.  <b>Question:</b> Does this imply that Liam didn't finish his work early?
BD	p(x): someone drinks lots of water q(x): they will feel hydrated r(x): they eat too much sugar s(x): they will experience a sugar crash p(a): Jane drinks lots of water ~p(a): Jane does not drink lots of water q(a): she will feel hydrated ~q(a): she will not feel hydrated r(a): she eats too much sugar ~r(a): she does not eat too much sugar s(a): she will experience a sugar crash ~s(a): she will not experience a sugar crash	<b>Context:</b> If an individual consumes a significant amount of water, they will experience a state of hydration. Conversely, if excessive amounts of sugar are ingested, a sugar crash will ensue. It is known that at least one of the following statements is true: either the Jane consumes ample water or she will not experience a sugar crash. However, the actual veracity of either statement remains ambiguous, as it could be the case that only the first statement is true, only the second statement is true, or both statements are true.  <b>Question:</b> Can we say at least one of the following must always be true? (a) she will feel hydrated and (b) she doesn't eat too much sugar
MP	p(x): someone is exhausted. q(x): they will take a rest. p(a): Jack is exhausted. ~p(a): Jack is not exhausted. q(a): he will take a rest. ~q(a): he will not take a rest.	<b>Context:</b> If someone is extremely tired, then they will seek some rest and relaxation. Today, Jack finds himself utterly exhausted.  <b>Question:</b> Does this entail that he will take rest?
DS	p: Levi is not studying for his exams ~p: Levi is studying for his exams q: Maya is writing a book ~q: Maya is not writing a book	<b>Context:</b> Either Levi is not studying for his exams or Maya is writing a book, or maybe both. It was unclear which choice to make as he didn't know if either of the options was true. However, it turned out that Levi decided to prioritize his exams and focus on studying. He knew that in order to succeed, he needed to dedicate his time and energy to preparing for the upcoming tests.  <b>Question:</b> Does this mean that Maya is writing a book?
HS	p(x): someone buys all the necessary supplies. q(x): they can start the project. r(x): they will finish it on time. p(a): Lily bought all the necessary supplies. ~p(a): Lily did not buy all the necessary supplies. q(a): She can start the project. ~q(a): She can not start the project. s(a): She will finish it on time. ~s(a): She will not finish it on time.	<b>Context:</b> If all the necessary supplies have been purchased by someone, then they can initiate the project. Once the project is started by someone, they will complete it within the expected timeframe.  <b>Question:</b> If Lily bought all the necessary supplies, does this mean that she will finish it on time?
CD	p: Harry goes to the park. ~p: Harry does not go to the park. q: he will have a picnic with his family ~q: he will not have a picnic with his family r: he goes to the beach ~r: he does not go to the beach s: he will swim in the ocean ~s: he will not swim in the ocean	<b>Context:</b> Harry had two options for his day off: going to the park or going to the beach. If he decided to go to the park, it meant that he would have a lovely picnic with his family. On the other hand, if he chose the beach, he would be able to enjoy a refreshing swim in the ocean. The only thing that was certain was that either Harry would go to the park or he would go to the beach. It was unclear which option he would ultimately choose, as it could be that he would go to the park alone, go to the beach alone, or even do both activities.  <b>Question:</b> Can we say at least one of the following must always be true? (a) he will have a picnic with his family and (b) he will swim in the ocean
DD	p: I order takeout ~p: I did not order takeout q: I will save time ~q: I will not save time r: I cook a meal ~r: I did not cook a meal s: I will save money ~s: I will not save money	<b>Context:</b> If I decide to order takeout, it means I will save time. On the other hand, if I choose to cook a meal, it means I will save money. The interesting thing is that I am in a situation where I won't be able to save time or money. It is uncertain whether I won't save time or I won't save money, or it could even be both scenarios. The only thing that is clear is that at least one of these possibilities is true.  <b>Question:</b> Can we say at least one of the following must always be true? (a) I don't order takeout and (b) I don't cook a meal
CT	p: Tom is an avid reader ~p: Tom is not an avid reader q: he devours books of all genres ~q: he does not devour books of all genres	<b>Context:</b> At least one of two things is true about Tom - he is either an avid reader or he devours books of all genres. We are unsure which one of these statements is true or if both are true. It could be that only the first statement is true, or only the second statement is true, or even that both are true.  <b>Question:</b> Can we say at least one of the following must always be true? (a) he devours books of all genres or (b) Tom is an avid reader
MI	p: he is not eating healthy ~p: he is eating healthy q: he will not gain weight ~q: he will gain weight	<b>Context:</b> Once upon a time, there was a guy named John who was very conscious about his weight. He knew that if he didn't eat healthy, he wouldn't be able to maintain his desired weight. John understood the condition that if he didn't eat healthy, he wouldn't be able to gain weight.  <b>Question:</b> Based on context, can we say, at least one of the following must always be true? (a) he is eating healthy and (b) he will not gain weight
EG	p(x): someone has coding skills a: Sheila p(a): Sheila is a proficient programmer	<b>Context:</b> Sheila possesses exceptional coding skills and is highly proficient in programming. She has mastered various programming languages and is known for her ability to write efficient and error-free code. Sheila's expertise in programming allows her to solve complex problems with ease. Her colleagues often seek her assistance when facing technical challenges as she consistently delivers high-quality solutions.  <b>Question:</b> does this mean that someone has coding skills?
UI	p(x): students need to take an exam to complete their degree a: Reema p(a): Reema is a student.	<b>Context:</b> All students are required to take an examination in order to fulfill the requirements for their degree. Reema, being a student, is also expected to fulfill the requirements.  <b>Question:</b> Does Reema need to take an exam to complete her degree?

Table 13: Illustrative examples of NL logical context and questions created using sentences that are generated in §3.2.1 for rules covered in PL and FOL for BQA. 13703

Rule	Generate Sentences in Step 1	NL logical expressions
DRI	<p>p: Cats and dogs are mammals.  q: Mammals typically have fur.  r: Cats don't have fur.  s: Dogs are loyal animals.  t: Dogs have fur.  ~t: Dogs don't have fur.</p>	<p><b>Context:</b> Once upon a time, in a land filled with animals, there were two popular mammalian creatures, cats and dogs. Mammals typically possessed a coat of fur, which kept them warm and protected. However, cats were an exception to this rule, as their bodies lacked fur. Nonetheless, both cats and dogs were beloved by many for their unique traits. Dogs, known for their loyalty, were particularly cherished by humans.</p> <p><b>Question:</b> Does this imply that dogs have fur?</p>
DRS	<p>p: John and Mary are parents.  q: Parents are usually loving and supportive.  r: Parents are normally responsible.  s: Mary isn't loving and supportive.  t: John is not responsible.  u: Mary is responsible.  ~u: Mary isn't responsible.  v: John is loving and supportive.  ~v: John isn't loving and supportive.</p>	<p><b>Context:</b> John and Mary were expecting their first child, filled with the anticipation and excitement that all parents feel. Parents are usually loving and supportive. Parents are normally responsible. However, something seemed amiss in their relationship. Mary, usually affectionate and caring, seemed distant and uninvolved. On the other hand, John, known for his responsible nature, started neglecting his duties and became unreliable.</p> <p><b>Question:</b> Does this imply that Mary is responsible and John is loving and supportive?</p>
DRD	<p>p: Jenny and Anna are tall.  q: Tall people usually play basketball.  r: Anna is possibly an exception to this rule.  s: Jenny plays basketball.  ~s: Jenny doesn't play basketball.</p>	<p><b>Context:</b> Jenny and Anna are known for their tall stature, which is often associated with playing basketball. However, Anna might be an exception to this norm.</p> <p><b>Question:</b> Does this entail that Jenny plays basketball?</p>
DRO	<p>p: Hummingbirds are birds.  q: Birds migrate south for the winter.  r: Hummingbirds do not migrate south for the winter.  s: All other birds than hummingbirds migrate south for the winter.  ~s: All other birds than hummingbirds don't migrate south for the winter.</p>	<p><b>Context:</b> In the bird kingdom, there are many different species that possess unique characteristics. One such species is the hummingbird, known for its ability to hover in mid-air and its vibrant colors. While most birds engage in the annual migration south for the winter, the hummingbird chooses to stay put and brave the cold weather. This decision sets the hummingbird apart from its fellow avian companions, as it relies on its resilience and resourcefulness to survive the harsh conditions.</p> <p><b>Question:</b> Does this mean that all other birds than hummingbirds migrate south for the winter?</p>
RE1	<p>p: Cats, dogs, and horses are animals.  q: Animals are usually considered to be intelligent creatures.  r: At least one of the cats or dogs is not considered intelligent.  s: Horses are considered to be intelligent creatures.  ~s: Horses aren't considered to be intelligent creatures.  t: Exactly one of the cats or dogs is not considered intelligent.  ~t: Exactly one of the cats or dogs is considered intelligent.</p>	<p><b>Context:</b> In a world where animals are often regarded as intelligent creatures, there is a captivating tale that revolves around cats, dogs, and horses. It is commonly believed that most animals possess a level of intellect. However, there is an intriguing twist to this belief as it is known that either cats or dogs are not considered particularly intelligent. As the story unfolds, we delve into the lives of these remarkable creatures, their interactions, and the unique qualities that each of them possesses.</p> <p><b>Question:</b> Does this entail that horses are considered to be intelligent creatures and exactly one of the cats or dogs is not considered intelligent?</p>
RE2	<p>p: cats normally meow.  q: At least one species of cat doesn't meow.  r: Exactly one species of cat doesn't meow.  ~r: Exactly one species of cat meows.</p>	<p><b>Context:</b> In the realm of cat communication, meowing serves as a fundamental aspect of their vocal repertoire. However, intriguingly enough, there exists a distinct species of cat that deviates from this conventional norm. This peculiar feline defies the expectations associated with its kind by refraining from emitting any meows whatsoever.</p> <p><b>Question:</b> Does this entail that exactly one species of cat doesn't meow?</p>
RE3	<p>p: Cars have four wheels.  q: wheels normally have spokes.  r: at least one wheel does not have spokes.  s: Cars have four wheels with spokes.  ~s: Cars don't have four wheels with spokes.</p>	<p><b>Context:</b> In a world where cars were known for having four wheels, it was considered a common fact that wheels typically came equipped with spokes. However, amidst this widespread understanding, there was an exception. At least one wheel defied this norm and stood out from the rest by not having any spokes at all.</p> <p><b>Question:</b> Does this imply that cars have four wheels with spokes?</p>
RAP	<p>p: John asserts that Sally was in the store.  q: Jane asserts that Sally was not in the store.  r: John's evidence is more reliable than Jane's.  ~r: John's evidence is less reliable than Jane's.  s: Sally was in the store.  ~s: Sally wasn't in the store.</p>	<p><b>Context:</b> In the midst of a heated argument, John adamantly claims that Sally was present at the store. However, Jane strongly opposes John's assertion, insisting that Sally was indeed absent from the store.</p> <p><b>Question:</b> If John's evidence is more reliable than Jane's, does this mean that Sally was in the store?</p>

Table 14: Illustrative examples of NL logical context and questions created using sentences that are generated in §3.2.1 for NM logic for BQA.

Type	Rule	Llama-2		Mistral		Gemini		ChatGPT		GPT-4	
		$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$
PL	HS	100.0	25.3	74.3	16.7	83.3	26.9	84.2	27.9	93.8	53.1
	DS	66.7	26.9	76.6	66.7	95.2	32.8	87.5	33.3	92.2	93.8
	CD	100.0	55.6	97.1	42.2	98.1	73.1	100.0	57.1	100.0	80.0
	DD	93.8	79.2	61.1	59.1	94.7	90.0	64.3	83.3	64.3	83.3
	BD	76.9	25.4	76.4	37.5	77.8	24.6	82.4	30.4	85.1	76.9
	MT	100.0	25.6	89.3	32.7	92.9	28.1	100.0	30.3	91.5	48.5
	MI	55.6	51.6	57.9	57.1	66.7	66.7	64.0	73.3	90.5	94.7
	CT	84.7	52.6	97.4	45.2	97.4	48.7	85.7	36.8	96.7	90.0
	<b>Avg</b>	<b>84.7</b>	<b>42.8</b>	<b>78.8</b>	<b>44.7</b>	<b>88.3</b>	<b>48.9</b>	<b>83.5</b>	<b>46.6</b>	<b>89.2</b>	<b>77.5</b>
FOL	EG	90.9	27.5	76.8	36.4	88.9	27.6	77.4	29.6	90.4	53.6
	UI	75.0	25.8	79.7	83.3	88.0	31.5	85.7	33.3	92.1	88.2
	MP	100.0	69.0	100.0	64.5	95.2	100.0	100.0	100.0	100.0	95.2
	HS	100.0	100.0	100.0	83.3	100.0	100.0	100.0	100.0	100.0	100.0
	DS	100.0	51.3	97.3	44.2	87.5	65.0	100.0	50.0	92.9	66.7
	CD	90.5	94.7	61.9	63.2	81.0	84.2	81.0	84.2	74.1	100.0
	DD	80.0	25.3	74.7	0.0	88.0	30.9	78.3	29.4	85.1	76.9
	BD	72.7	58.6	77.8	58.1	81.8	94.1	79.2	93.8	90.5	94.7
	<b>Avg</b>	<b>88.3</b>	<b>57.9</b>	<b>80.7</b>	<b>55.2</b>	<b>88.1</b>	<b>67.5</b>	<b>87.5</b>	<b>68.3</b>	<b>88.0</b>	<b>86.2</b>
NM	DRI	86.2	70.6	62.5	62.5	80.4	90.9	74.0	90.0	95.2	100.0
	DRS	47.4	47.6	77.3	83.3	61.9	70.6	60.7	75.0	90.0	90.0
	DRD	100.0	66.7	100.0	62.5	78.6	68.0	63.6	55.2	70.4	92.3
	DRO	58.1	14.3	66.7	7.7	62.8	9.1	64.1	14.6	76.5	33.3
	RE1	60.9	64.7	58.6	72.7	66.7	85.7	60.6	100.0	87.0	100.0
	RE2	0.0	25.0	88.0	30.9	84.8	32.6	77.8	25.8	86.5	46.4
	RE3	52.6	52.4	68.0	80.0	82.4	71.4	62.1	81.8	61.5	71.4
	RAP	50.0	50.0	65.5	90.9	62.5	100.0	55.6	100.0	87.0	100.0
	<b>Avg</b>	<b>56.9</b>	<b>48.9</b>	<b>73.3</b>	<b>61.3</b>	<b>72.5</b>	<b>66.0</b>	<b>64.8</b>	<b>67.8</b>	<b>81.8</b>	<b>79.2</b>

Table 15: Performance of LLMs in few-shot setting in terms of label-wise accuracy on LogicBench(Eval), where  $A(Yes)$  and  $A(No)$  denote the accuracy for the *Yes* and *No* labels, respectively.

Rule	Context/Question	Actual/Pred	Reasoning Chains
PL(MT)	Context: if elijah has a new job, then he will get a raise. however, he would not be receiving the raise he had anticipated. Question: does this infer that elijah doesn't have a new job?	Yes/No	<ol style="list-style-type: none"> <li>1. The context provides a conditional statement: If Elijah has a new job, then he will get a raise.</li> <li>2. The context also states that Elijah would not be receiving the raise he had anticipated.</li> <li>3. The second statement does not negate the first; it only states that the raise is not as expected.</li> <li>4. The lack of the anticipated raise does not necessarily mean Elijah doesn't have a new job; it could mean that the raise with the new job is different from what he expected.</li> <li>5. Therefore, the information provided does not allow us to infer that Elijah doesn't have a new job based solely on the raise situation.</li> </ol>
PL(MT)	Context: james had always been forgetful, especially when it came to wearing a sweater. his mother had repeatedly told him, 'if he forgets to wear a sweater, then he will catch a cold.' throughout the day, he went about his activities, oblivious to the dropping temperature. hours later, he returned home, but james didn't catch a cold. Question: does this infer that james didn't forget to wear a sweater?	Yes/No	<ol style="list-style-type: none"> <li>1. The mother's rule establishes a conditional relationship: If James forgets to wear a sweater (condition), then he will catch a cold (result).</li> <li>2. It is stated that James did not catch a cold upon returning home.</li> <li>3. According to the mother's rule, if the result (catching a cold) did not occur, we cannot conclusively determine whether the condition (forgetting to wear a sweater) happened or not.</li> <li>4. The absence of the result (catching a cold) does not necessarily mean the condition (forgetting to wear a sweater) did not occur, as there could be other factors at play that prevented James from catching a cold.</li> <li>5. Therefore, from the information given, we cannot infer with certainty that James remembered to wear a sweater.</li> </ol>
FOL(DD)	Context: if an individual operates a vehicle with caution, they will prevent any accidents from occurring. however, if they decide to accelerate and exceed the speed limit, they will be issued a fine. there is a certainty that at least one of the following statements is true: either john will not be able to avoid an accident or he will not face a penalty. it is unclear which of these statements is accurate, as it could be only one of them, or it is possible that both of them hold true. Question: can we say at least one of the following must always be true? (a) john doesn't drive carefully and (b) he doesn't speed up	Yes/No	<ol style="list-style-type: none"> <li>1. According to the context, if an individual operates a vehicle with caution, they will prevent any accidents from occurring.</li> <li>2. Also, if they decide to accelerate and exceed the speed limit, they will be issued a fine.</li> <li>3. It is stated that there is a certainty that at least one of the following statements is true: either John will not be able to avoid an accident or he will not face a penalty.</li> <li>4. It is unclear which of these statements is accurate, as it could be only one of them, or it is possible that both of them hold true.</li> <li>5. The question asks if we can say at least one of the following must always be true: (a) John doesn't drive carefully and (b) he doesn't speed up.</li> <li>6. Based on the given information, we cannot say that at least one of the statements (a) and (b) must always be true.</li> </ol>
FOL(DD)	Context: if someone is taking care of thier health, then they will be fit. however, indulging in unhealthy habits can make individuals susceptible to various diseases. the truth is, we can be certain about at least one of the following possibilities: either jenny won't be fit and healthy, or she won't be prone to diseases. it is important to note that we are unaware of which statement is accurate. it could be the case that only the first statement is true, only the second statement is true, or both statements hold validity. Question: can we say at least one of the following must always be true? (a) jenny takes care of her health and (b) she doesn't indulge in unhealthy habits	No/Yes	<ol style="list-style-type: none"> <li>1. According to the context, if someone takes care of their health, they will be fit.</li> <li>2. Indulging in unhealthy habits can make individuals susceptible to various diseases.</li> <li>3. The context states that we can be certain about at least one of the following possibilities: either Jenny won't be fit and healthy or she won't be prone to diseases.</li> <li>4. It is important to note that we are unaware of which statement is accurate.</li> <li>5. The possibilities are: <ol style="list-style-type: none"> <li>a) Jenny takes care of her health and is fit.</li> <li>b) Jenny doesn't indulge in unhealthy habits and is not prone to diseases.</li> <li>c) Both statements hold validity.</li> </ol> </li> <li>6. Therefore, at least one of the following must always be true: (a) Jenny takes care of her health and (b) she doesn't indulge in unhealthy habits.</li> </ol>

Table 16: Case study of reasoning chains for PL(MT) generated from ChatGPT, and FOL(DD) generated from GPT-4.

Logic Type	Rule	BQA		MCQA
		A(No)	A(Yes)	Acc.
PL	HS	98.81 <sub>0.01</sub>	92.13 <sub>0.02</sub>	96.67 <sub>0.03</sub>
	DS	63.61 <sub>0.06</sub>	87.25 <sub>0.07</sub>	90.00 <sub>0.05</sub>
	CD	88.88 <sub>0.02</sub>	41.04 <sub>0.05</sub>	86.67 <sub>0.03</sub>
	DD	73.69 <sub>0.01</sub>	19.53 <sub>0.01</sub>	63.33 <sub>0.08</sub>
	BD	83.62 <sub>0.03</sub>	43.38 <sub>0.04</sub>	65.00 <sub>0.05</sub>
	MT	46.88 <sub>0.03</sub>	31.55 <sub>0.05</sub>	58.33 <sub>0.06</sub>
	MI	84.49 <sub>0.01</sub>	31.57 <sub>0.01</sub>	71.67 <sub>0.06</sub>
	CT	95.48 <sub>0.04</sub>	59.15 <sub>0.01</sub>	70.00 <sub>0.13</sub>
	<b>Avg</b>	79.43 <sub>0.03</sub>	50.7 <sub>0.03</sub>	75.21 <sub>0.06</sub>
FOL	EG	100.0 <sub>0.0</sub>	98.33 <sub>0.03</sub>	88.33 <sub>0.03</sub>
	UI	90.32 <sub>0.00</sub>	94.64 <sub>0.00</sub>	88.33 <sub>0.08</sub>
	MP	84.87 <sub>0.02</sub>	100.0 <sub>0.0</sub>	93.33 <sub>0.08</sub>
	HS	96.67 <sub>0.02</sub>	94.52 <sub>0.00</sub>	76.67 <sub>0.03</sub>
	DS	77.42 <sub>0.03</sub>	90.28 <sub>0.09</sub>	55.00 <sub>0.05</sub>
	CD	90.41 <sub>0.01</sub>	40.39 <sub>0.02</sub>	81.67 <sub>0.08</sub>
	DD	72.77 <sub>0.02</sub>	22.20 <sub>0.03</sub>	50.00 <sub>0.20</sub>
	BD	79.60 <sub>0.03</sub>	30.05 <sub>0.06</sub>	60.00 <sub>0.05</sub>
	<b>Avg</b>	82.56 <sub>0.02</sub>	69.61 <sub>0.05</sub>	73.52 <sub>0.09</sub>
NM	DRI	86.10 <sub>0.08</sub>	98.25 <sub>0.03</sub>	60.00 <sub>0.09</sub>
	DRS	73.79 <sub>0.00</sub>	10.37 <sub>0.10</sub>	30.00 <sub>0.05</sub>
	DRD	85.76 <sub>0.05</sub>	100.0 <sub>0.0</sub>	55.00 <sub>0.09</sub>
	DRO	70.37 <sub>0.04</sub>	100.0 <sub>0.0</sub>	33.33 <sub>0.06</sub>
	REI	83.80 <sub>0.03</sub>	44.31 <sub>0.09</sub>	73.33 <sub>0.08</sub>
	REII	63.27 <sub>0.02</sub>	76.03 <sub>0.06</sub>	63.33 <sub>0.08</sub>
	REIII	63.89 <sub>0.03</sub>	88.33 <sub>0.13</sub>	56.67 <sub>0.03</sub>
	RAP	68.15 <sub>0.06</sub>	87.12 <sub>0.03</sub>	38.33 <sub>0.08</sub>
	<b>Avg</b>	74.39 <sub>0.04</sub>	75.55 <sub>0.05</sub>	51.25 <sub>0.07</sub>

Table 17: Evaluation of Yi-34B in terms of accuracy on  $LogicBench(Eval)_{BQA}$  and  $LogicBench(Eval)_{MCQA}$ .