
LogicBench: A Benchmark for Evaluation of Logical Reasoning

Anonymous Author(s)

Abstract

1 Recently developed large language models (LLMs) have been shown to perform
2 remarkably well on a wide range of language understanding tasks. But, can they
3 really “Reason” over the natural language? This question has been receiving signif-
4 icant research attention and a number of reasoning skills such as commonsense,
5 numerical, and qualitative have been studied. However, the crucial skill pertaining
6 to ‘logical reasoning’ has remained underexplored. Existing work investigating
7 this reasoning ability has focused only on a couple of axioms (such as modus
8 ponens and modus tollens) of propositional and first-order logic. To study logical
9 reasoning, we introduce *LogicBench*, a systematically created natural language
10 question-answering dataset encompassing 25 reasoning patterns spanning over
11 propositional, first-order, and non-monotonic logics. Key steps of our dataset
12 construction consist of (1) controlled generation of sentences and their negations
13 containing different ontologies, (2) (*context, question, answer*) triplets creation us-
14 ing heuristically designed templates, and (3) semantic variations of triplets adding
15 more diversity. We first evaluate easily accessible and widely used LLMs such as
16 GPT-3, ChatGPT, and FLAN-T5 and show that they do not fare well on *LogicBench*,
17 achieving just above random accuracy on average ($\sim 52\%$). Then, we show that
18 LLMs trained using our data exhibit a better understanding of logical reasoning
19 leading to performance improvements on several existing logical reasoning datasets
20 such as LogiNLI, FOLIO, LogiQA, and ReClor.¹

21 1 Introduction

22 Large language models such as GPT-3 [3], ChatGPT, and FLAN [18] have made remarkable progress
23 in NLP research enabling machines to perform a variety of language tasks that were previously
24 thought to be exclusive to humans [12, 2, 20]. However, the ability of these LLMs to reason
25 “logically” over natural language text remains under-explored, even though logical reasoning is a
26 fundamental aspect of intelligence and a crucial requirement for many practical applications, such
27 as question-answering systems [8] and conversational agents [1]. Although several datasets have
28 been proposed [4, 16, 7, 13] to evaluate the logical reasoning capabilities of LLMs, these datasets
29 are limited in their scope by (1) not evaluating logical reasoning independently of other forms of
30 reasoning such as LogiQA [11] and ReClor [19]; and (2) evaluating only a single type of logic and
31 covering only few logical inference rules as done in FOLIO [6] and ProntoQA [14]. Thus, our aim in
32 this work is to address the lacuna of having a more comprehensive evaluation dataset for LLMs.

33 To this end, we propose *LogicBench*, a systematically created question-answering dataset for the
34 evaluation of logical reasoning ability. As illustrated in Figure 1, *LogicBench* includes a total of 25

¹Data is available at <https://anonymous.4open.science/r/LogicBench-EEBB>

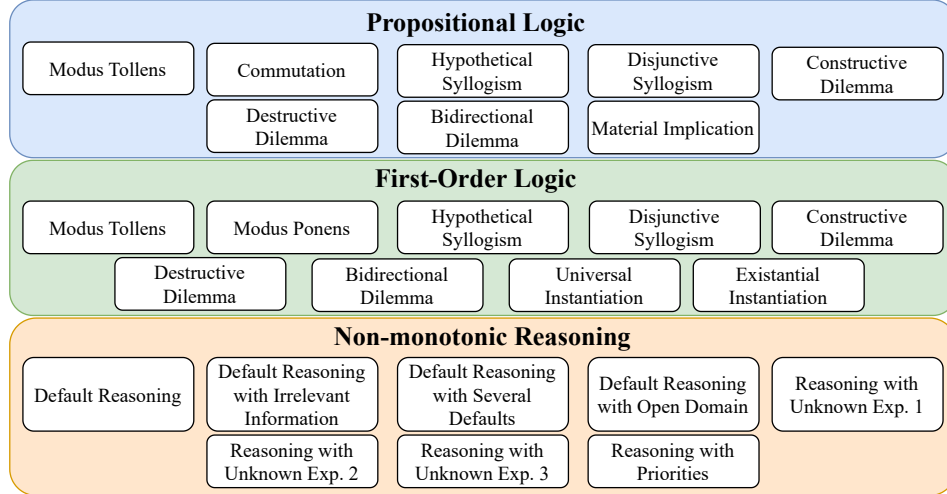


Figure 1: Comprehensive representation of different inference rules and reasoning patterns covered by propositional, first-order, and non-monotonic logics. *Exp.* indicates Expectation

35 reasoning patterns across propositional, first-order, and non-monotonic logics. To evaluate LLMs, we
 36 formulate a binary classification task in *LogicBench* in which the context represents logical statements
 37 and the models have to determine whether a conclusion given in the question is logically entailed by
 38 the context. For example, given the context “All mammals have fur” and “A cat is a mammal”, for
 39 the question is “Does a cat have fur?”, the correct answer, is "Yes". (Additional examples of task
 40 instances are presented in Table 3 and Appendix B. To construct *LogicBench*, we use a three-stage
 41 procedure (refer to §2). In the first stage, we prompt GPT-3 to generate a variety of coherent natural
 42 language sentences consisting of different ‘ontologies’ (i.e., a collection of concepts such as car,
 43 person, and animals) and their corresponding negations (refer to §2.2.1). Then, in the second stage,
 44 we generate (*context, question, answer*) triplets using heuristically designed templates based on
 45 the inference rules and patterns. Finally, in the third stage, we generate semantics preserving and
 46 inverting variations of these logical rules by incorporating negations.

47 We evaluate a range of accessible and widely used LLMs including GPT-3 [3], ChatGPT, FLAN-T5
 48 [18], Tk-instruct [17], and UnifiedQA [9] with respect to *LogicBench* on the accuracy of the predicted
 49 answers (i.e., “Yes” or “No”). Experimental results reveal that these models struggle with respect
 50 to many of the inference rules and patterns (showing $\sim 52\%$ accuracy on an average), suggesting
 51 significant room for improvement in their logical reasoning abilities. We then synthetically augment
 52 *LogicBench* and train T5-large. Our initial experimental results show that this improves the logical
 53 reasoning ability of existing models leading to performance improvement on other logic datasets, such
 54 as LogicNLI, and FOLIO ($\sim 2\%$ on an average), and shows competitive performance on LogiQA
 55 and ReClor. In summary, our contributions are as follows:

- 56 1. Introducing *LogicBench*: A systematically created dataset to assess the logical reasoning
 57 capabilities of LLMs across propositional, first-order, and non-monotonic logics. This
 58 benchmark will be publicly available for evaluation and training purposes.
- 59 2. We propose a three-stage method to construct *LogicBench* consisting of GPT-3 to generate
 60 coherent natural language sentences using prompts and a template-based module to convert
 61 them into logical rules. By assessing the performance of existing LLMs, we gain insights
 62 into their logical reasoning abilities which further leads to several interesting findings.
- 63 3. To the best of the authors’ knowledge, this is the first benchmark to study non-monotonic
 64 reasoning, as well as various inference rules in propositional and first-order logics including
 65 hypothetical and disjunctive syllogism; and bidirectional, constructive, and destructive
 66 dilemmas in the NLP domain.

67 2 LogicBench

68 In this section we discuss the logic types, inference rules, and patterns that are explored in this
69 research. We also outline the methods for generating the data, and statistics of *LogicBench*.

70 2.1 Logics Types

71 **Propositional Logic (PL)** Propositional logic employs a collection of statements or propositions
72 (denoted as $\mathcal{P} = p_1, p_2, \dots, p_n$, where p_i represents a proposition) and builds upon them using logical
73 connectives such as ‘ \wedge ’, ‘ \vee ’, ‘ \rightarrow ’, ‘ \leftrightarrow ’, and ‘ \neg ’. Several inference rules for propositional logic
74 have been defined using which given a set of premises, one can derive a sound conclusion. To
75 illustrate this, let us consider two propositions: p_1 , which states "It is raining," and p_2 , which states
76 "It is cloudy." From these propositions, we can construct a context (KB) consisting of two premises:
77 (1) $p_1 \rightarrow p_2$ and (2) p_1 . Based on this KB, we can conclude p_2 . This inference rule is written as
78 $((p_1 \rightarrow p_2) \wedge p_1) \vdash p_2$ and is known as ‘Modus Ponens’. In our study, we explore nine distinct
79 inference rules of propositional logic, extensions of seven of them with one-variable and a universal
80 quantifier, and two axioms of first-order logic as shown in Table 1. These inference rules provide a
81 systematic framework for deriving valid conclusions.

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	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (p \vee r)) \vdash (q \vee s)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (p(a) \vee r(a)))) \vdash (q(a) \vee s(a))$
DD	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (\neg q \vee \neg s)) \vdash (\neg p \vee \neg r)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (\neg q(a) \vee \neg s(a)))) \vdash (\neg p(a) \vee \neg r(a))$
BD	$((p \rightarrow q) \wedge (r \rightarrow s) \wedge (p \vee \neg s)) \vdash (q \vee \neg r)$	$(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x)) \wedge (p(a) \vee \neg s(a)))) \vdash (q(a) \vee \neg r(a))$
CT	$(p \vee q) \vdash (q \vee p)$	-
MI	$(p \rightarrow q) \vdash (\neg p \vee q)$	-
EI	-	$\exists x P(x) \Rightarrow P(a)$
UI	-	$\forall x A \Rightarrow A\{x \mapsto a\}$

Table 1: Inference rules and (two) axioms that establish the relationship between premises and their corresponding 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, EI: Existential Instantiation, UI: Universal Instantiation

82 **First-order Logic (FOL)** In this work, we consider a restricted set of logical axioms for FOL that
83 utilize quantifiers, \forall (universal quantifier) and \exists (existential quantifier). The universal quantifier
84 (\forall) denotes that a statement holds true for all instances within a specific category. In contrast, the
85 existential quantifier (\exists) indicates that a statement is true for at least one instance within its scope.
86 For instance, a simple extension of propositional ‘Modus Ponens’ is an inference rule where given
87 the premises $\forall x(p(x) \rightarrow q(x))$ and $p(a)$, we conclude $q(a)$ (e.g., given “All kings are greedy” and
88 “Sam is a king”, we can conclude “Sam is greedy”). Here, we explore various axioms and inference
89 rules that incorporate the quantifiers shown in Table 1.

90 **Non-monotonic (NM) Reasoning** In this work, we analyze a range of logical reasoning templates
91 in NM logics involving “Default Reasoning,” “Reasoning about Unknown Expectations,” and “Reasoning
92 about Priorities.” These templates are inspired by the compilation [10] made in 1989 to
93 evaluate the abilities of various non-monotonic logics that were being developed at that time. Below
94 Table 2 shows examples of NM reasoning. Additional examples are given in Appendix B.3.

95 A key aspect of NM logics is to formalize notions such as "normally," "typically," and "usually"
96 that are not directly formalizable using classical quantifiers in the first-order setting. The general
97 rule “Heavy blocks are normally located on the table” does not imply that “All heavy blocks are

Basic Default Reasoning	Default Reasoning with Irrelevant Information
Context: Blocks A and B are heavy. Heavy blocks are typically located on the table. A is not on the table. Conclusion: B is on the table.	Context: Blocks A and B are heavy. Heavy blocks are typically located on the table. A is not on the table. B is red. Conclusion: B is on the table.
Reasoning about Unknown Expectations	Reasoning about Priorities
Context: Blocks A, B, and C are heavy. Heavy blocks are normally located on the table. At least one of A, B is not on the table. Conclusion: C is on the table. Exactly one of A, B is not on the table.	Context: Jack asserts that block A is on the table. Mary asserts that block A is not on the table. When people assert something, they are normally right. Conclusion: If Mary's evidence is more reliable than Jack's, then block A is not on the table

Table 2: Illustrative examples of non-monotonic reasoning adapted from [10]

98 always located on the table". Rather, this rule allows for exceptions. Our work explores various NM
 99 reasoning types, as depicted in Figure 1, to delve deeper into the nuances of this type of reasoning.

100 2.2 Data Creation

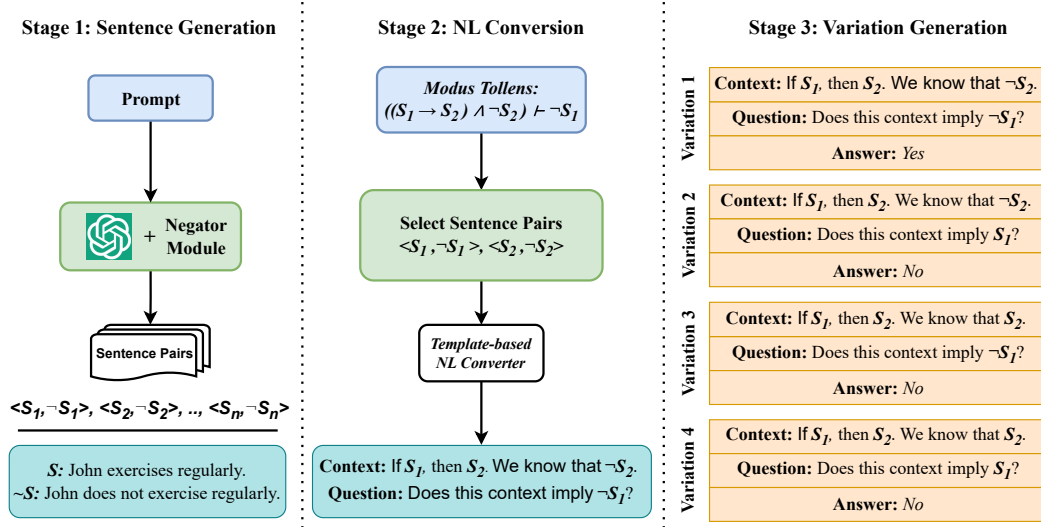


Figure 2: Schematic representation of three-stage procedure for data creation. NL: Natural Language

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

- 102 1. **Sentence Generation:** Starting with a given prompt, we generate coherent sentences and
 103 their negations that incorporate different ontologies.
- 104 2. **NL Conversion:** Using predefined templates of reasoning patterns based on their formal
 105 expressions, we convert the generated sentences into *(context, question, answer)* triplets.
- 106 3. **Variation Generation:** We generate semantically preserving and inverting variations of
 107 these triplets to add more diversity.

108 By following this method, we construct *LogicBench*, and examples of generated data corresponding
 109 to each logic type and reasoning patterns are presented in Appendix B.

110 2.2.1 Sentence Generation

111 Here, the first step is to generate sentences with diverse *ontologies*. An ontology represents a
 112 collection of concepts (e.g. car, person, animals, etc.) along with their corresponding associated

113 properties. To generate these sentences, we prompt the GPT-3 model with instructions tailored for
 114 each inference rule. The prompt schema, as depicted in Figure 3, comprise three crucial components:

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

118 **Examples** provide sample sentences that need to be gener-
 119 ated. We also illustrate how these sentences will be utilized
 120 in later stages, emphasizing the importance of coherence
 121 and the inclusion of relevant ontological concepts.

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

124 An example of a prompt corresponding to the ‘Modus
 125 Tollens’ from PL is presented in Appendix A for better
 126 illustration. Note that our objective at this stage is not to
 127 generate logical sentences but rather to generate a diverse
 128 and coherent set of sentences that encompass various con-
 129 cepts. We also create a negation sentence corresponding to
 130 each generated sentence². In this work, the scope of generating negations is simple (refer to Appendix
 131 C for examples), however, negations can be more complicated in the case of logic. These generated
 132 sentences will be combined with logical connectives in a later stage to form context and questions.

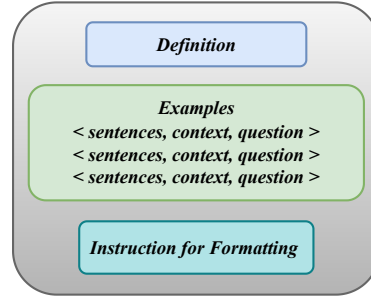


Figure 3: Schematic representation of prompt.

133 2.2.2 NL Conversion

134 We focus on leveraging the formal expressions of reasoning patterns to create templates that establish
 135 the desired NL formulation for each logical connective. For instance, implication: “ $p \rightarrow q$ ” is
 136 expressed as “If p , then q ”, conjunction: “ $p \wedge q$ ” is expressed as “ p and q .”, and disjunction: “ $p \vee q$ ”
 137 is expressed as “At least one of the following is true: (1) p and (2) q . Note that we do not know which
 138 of (1) and (2) is true. It is possible that only (1) is true, or only (2) is true, or both are true.”

139 With these established formulations, we proceed to utilize the sentences generated in §2.2.1 to
 140 create the context and questions corresponding to reasoning patterns. For instance, let’s consider
 141 the “Modus Tollens” from PL ($((p \rightarrow q) \wedge \neg q) \vdash \neg p$), and the “Bidirectional Dilemma” from FOL
 142 ($(\forall x((p(x) \rightarrow q(x)) \wedge (r(x) \rightarrow s(x))) \wedge (p(a) \vee \neg s(a))) \vdash (q(a) \vee \neg r(a)))$). Table 3 presents
 143 examples of logical context and questions for these inference rules, and Appendix C showcases
 144 further examples corresponding to each inference rule and patterns from *LogicBench*.

145 2.2.3 Variation Generation

146 After generating the context and questions in §2.2.2, we generate semantically preserving and
 147 inverting variations of questions. Let’s consider the example of “Modus Tollens” from Table 3,
 148 where the question is: “If he won’t order pizza for dinner, does this imply that Liam didn’t finish his
 149 work early?” In this question, we observe two propositions: s_1 , representing the statement “Liam
 150 didn’t finish his work early,” and s_2 , representing the statement “He won’t order pizza for dinner.”
 151 By perturbing these propositions, we can create four possible tuples: $\langle s_1, s_2 \rangle, \langle \neg s_1, s_2 \rangle$
 152 $\langle s_1, \neg s_2 \rangle, \langle \neg s_1, \neg s_2 \rangle$. Each tuple represents a combination of true or negation values
 153 for the propositions. Although it is possible to create more combinations from $\langle s_1, \neg s_1 \rangle$, and
 154 $\langle s_2, \neg s_2 \rangle$, we refine and restrict the set of triplets to exclude those that undermine the validity
 155 of the inference rule. To generate question variations, we replace the propositions in the original
 156 question with the corresponding tuples from the generated variations, hence, adding more diversity
 157 to *LogicBench*. This process allows us to create different variations of the question, as illustrated in
 158 Figure 2 (Step 3). More examples of question variations are in Appendix B.

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

Axiom	Generated sentences in stage 1	Context and Question
Modus Tollens	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.	Context: If Liam finished his work early, then he will order pizza for dinner. Question: If he won't order pizza for dinner, does this imply that Liam didn't finish his work early?
Bidirectional Dilemma	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	Context: 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. Question: If at least one of (1) and (2) is true, 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.

Table 3: Illustrative examples of logical context and questions created using sentences that are generated in the first stage §2.2.1.

159 2.3 Statistics and Qualitative Analysis

160 **Statistics** We introduce two versions of our proposed dataset: *LogicBench(Eval)* and *Log-*
161 *icBench(Aug)*. Statistics of both versions are presented in Table 4. Here, *LogicBench(Eval)* is
162 created using the above method along with human-in-loop to ensure the quality of generated data,
163 whereas *LogicBench(Aug)* is only a synthetically augmented version for training purposes.

164 These two versions aim
165 to accommodate differ-
166 ent evaluation and train-
167 ing needs to explore log-
168 ical reasoning. Consider-
169 ing the cost and complex-
170 ity associated with recent
171 LLMs such as GPT-3, and
172 GPT-4, we believe that *LogicBench(Eval)* provides a more feasible evaluation benchmark.

Dataset	# of Instances per Axiom	Total # of Instances	Total # of Instances (Including Variations)
<i>LogicBench(Eval)</i>	20	500	1720
<i>LogicBench(Aug)</i>	150	3750	12908

Table 4: Statistics of the *LogicBench(Eval)* and *LogicBench(Aug)*

173 **Quality of Data** Throughout the data generation phase of *LogicBench(Eval)*, the authors conduct
174 a review of the logical formations to ensure they adhered to the intended structure. We examine
175 each reasoning pattern for any potential discrepancies, ensuring that they were logically sound and
176 correctly represented the intended relationships between propositions. In addition to the logical
177 formation, we also dedicated considerable effort to eliminating typos and validating the grammar.

178 3 Results and Analysis

179 3.1 Experimental Setup

180 **Task Formulation** We formulate binary classification task using *LogicBench* to evaluate the logical
181 reasoning ability of LLMs. Let us consider a set of data instances $\mathcal{I}_{a,L}$ corresponding to axiom a
182 and logic type L . In this set, i^{th} instance is represented as $\mathcal{I}_{a,L}^i = \{(c_i, Q_i)\}$ where c_i represents
183 context and $Q_i = \{q_1, q_2, \dots, q_n\}$ represents set of question and its variations corresponding to i^{th}
184 instance. As discussed in §2, each context (c) represents logical rules (e.g., All cats have fur. Tom is
185 a cat.) and question (q) represents the conclusion (e.g., Does Tom have fur?). To each context and
186 question pair, i.e., $\langle c, q \rangle$, we assign a label from the set $\mathcal{Y} = \{Yes, No\}$. We assign a label *Yes*
187 if the conclusion logically entails the context, otherwise, assign a label *No*. To evaluate any model
188 on this setup, we provide $\langle c, q \rangle$ as input to predict a label from \mathcal{Y} .

189 **Experiments** We evaluate easily available and widely used prompting models (i.e., GPT-3 (davinci-
190 003) and ChatGPT), and instruction-tuned models (FLAN-T5 and Tk-instruct) on *LogicBench(Eval)*.
191 Since logical reasoning is an important aspect of different QA tasks, we also evaluate UnifiedQA.
192 Each model is evaluated in a zero-shot setting where the prompt is provided to the model without
193 any in-context examples. This approach allows us to determine LLM’s inherent ability to do logical
194 reasoning (based on pre-training), as we can not expect that various logical inference rules/patterns
195 will always be made part of prompts. However, we do evaluate these models in a few-shot setting,
196 and present the results in Appendix F. We also present exploratory – only exploratory because of the
197 limited availability of their inference APIs – analysis over Bard and GPT-4 in Appendix G.

198 In addition, we employed the T5-large model and trained it on the *LogicBench(Aug)* resulting in a
199 model named LogicT5. LogicT5 has achieved $\sim 97\%$ of accuracy on *LogicBench(Eval)* since it is
200 evident that supervised fine-tuning improves results by a large margin. Subsequently, we performed
201 fine-tuning on four other logical reasoning datasets: LogiQA, Reclor, LogicNLI, and FOLIO. Our
202 experiments were carried out in two settings: single-task (fine-tuning and evaluation on one dataset)
203 and multi-task (fine-tuning on all four datasets combined, with separate evaluations for each dataset).
204 A detailed experimental setup is described in Appendix D.

205 **Metrics** Here, we evaluate performance in terms of accuracy corresponding to each label, i.e.,
206 $A(Yes)$ and $A(No)$. We evaluate each model on three different prompts and report average results
207 across these prompts. All prompts used for experiments are described in Appendix D.

208 3.2 Benchmark Results

209 Table 5 represents label-wise accuracy ($A(Yes)$ and $A(No)$) corresponding to each LLMs. Here,
210 we focus on analyzing the $A(Yes)$ since the aim is to understand the model’s logical reasoning
211 capabilities in answering the question where the conclusion entails the logical context. Table 5
212 provides valuable insights into the performance of different models on various logic types. For
213 PL, UnifiedQA exhibits an average performance of 15%, while FLAN-T5 and Tk-instruct achieve
214 $\sim 25\%$. GPT-3 demonstrates a performance of 57.6%, and ChatGPT achieves 46.8%. Moving on to
215 FOL, these models showcase performance accuracy of 52.7%, 51.2%, 55.7%, 76.2%, and 72.6% for
216 UnifiedQA, FLAN-T5, Tk-instruct, GPT-3, and ChatGPT, respectively. On the NM reasoning, these
217 models show an accuracy of 63.5%, 56.2%, 56.3%, 62%, and 70.9%, respectively. Overall, these
218 models display an average performance of $\sim 34\%$, $\sim 61\%$, and $\sim 62\%$ on PL, FOL, and NL.

219 From Table 5, we can observe that models struggle more with inference rules of PL compared to
220 FOL and NM reasoning. Furthermore, it is noticeable that each model performs relatively better on
221 questions with a negative response (i.e., No) compared to questions with a positive response (i.e.,
222 Yes). This observation suggests that the models struggle to fully comprehend the logical relationship
223 between the context and the conclusion (i.e., lower $A(Yes)$). However, they demonstrate a relatively
224 stronger understanding when the relationship is contradictory in nature (i.e., higher $A(No)$). However,
225 analyzing the performance of the models on inference rules is crucial to understand their limitations.
226 Table 5 presents the inference rule-wise performance for each model as well.

227 3.3 Analysis and Discussion

228 **Large models are better logical reasoners.** Based on the observed performance from Table 5,
229 it becomes evident that larger model sizes and extensive pre-training data contribute to a better
230 understanding of logical aspects. Consequently, models with larger sizes tend to exhibit higher
231 performance across different types of logic. Nonetheless, the average performance remains at around
232 52.7%, indicating room for improvement in these models’ logical comprehension capabilities.

233 **Negations are hard to understand when embedded with logical rules.** Regarding PL and FOL,
234 it is apparent that the models struggle more with the DS, DD, and MT inference rules. A closer
235 look at Table 1 reveals that all of these axioms include examples where the models need to draw
236 conclusions based on negated premises. This indicates that the models encounter difficulties when

Type	Axiom	FLAN-T5		Tk-instruct		UnifiedQA		GPT-3		ChatGPT	
		$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$	$A(No)$	$A(Yes)$
PL	HS	100	48.4	97.9	57.9	81.6	95.2	97.6	78.3	100	57.2
	DS	64.1	8.3	67.9	10.9	68.8	2.1	75.5	33.3	73.8	5.5
	CD	50	25	75	25	63.3	0	97.7	75.4	99.4	81.0
	DD	75	25	75	25	71.1	0	78	43.4	100	33.1
	BD	75	25	75	25	88.8	0	80.5	97	97.4	58.0
	MT	92.2	44.6	74.5	24.4	74.1	22.9	72.5	17.5	92.3	45.5
	MI	63.7	23.2	64.2	0	90.3	0	81.5	33.3	91.3	41.3
	CT	25	16.7	78.3	31.5	95.2	0	95.8	97	100	52.3
	Avg	68.1	27	76	25	79.1	15	84.9	59.4	94.3	46.8
FOL	EI	100	100	95	100	98.4	100	88.9	100	89.7	100
	UI	98.1	86.9	89.3	84.4	72.5	94.9	88.2	98.2	85.1	94.3
	MP	99.2	79.3	88.6	86.3	70.7	87.4	81.6	82.3	88.5	80.1
	HS	100	49.2	100	52.7	83.6	88.3	94.9	78.7	95.7	53.1
	DS	72.1	21.9	71.4	4.6	80.4	55.6	81.8	96.3	88.2	97.6
	CD	75	25	91.7	62	54.6	0	93.2	65.9	93.7	87.9
	DD	75	25	87.4	28	94.4	0	75.4	44.4	83.9	30.6
	BD	25	25	91.7	47	100	33.3	77.5	94.4	98.7	67.6
MT	93.3	48.1	81.8	35.9	70.8	15.2	74.8	25.7	85.9	42.3	
	Avg	82	51.2	88.5	55.7	80.6	52.7	84.1	76.2	89.9	72.6
NM	DRI	60.5	59.6	52.5	53.8	58.2	61.7	75	100	75.6	89.6
	DRS	66.3	2.9	60	3.9	67.3	2.8	72.6	10.1	72.7	0
	DRD	95	95	88.8	75.7	68.1	97.8	84.7	100	82.2	100
	DRO	40	42.6	43.8	45.3	53.2	91.7	65.3	100	70.3	100
	RE1	74.2	24.2	85.2	28	75.8	33.3	74.3	0	81.4	33.6
	RE2	100	100	98.2	93.8	56.2	66.7	50	0	62.3	64.7
	RE3	65.6	63	78.3	57.7	78.2	81	64.5	93.6	67.2	82.7
RAP	70.1	62.6	76.9	92.5	64.5	73	56.8	92.2	58.3	96.9	
	Avg	71.5	56.2	73	56.3	65.2	63.5	67.9	62	71.3	70.9

Table 5: Evaluation of LLMs 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. 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

237 negated premises are introduced. Additionally, the performance of the models tends to decrease when
238 inference rules involve negations.

239 **Longer inference rules are still challenging.** Table 1 indicates that the models face challenges
240 when handling longer rules, such as BD, CD, and DD, both in PL and FOL. Hence, it can be
241 concluded that these models struggle with longer logical dependencies in the premise, particularly
242 when a higher number of propositions are present. In the case of NM reasoning, the models exhibit
243 lower performance in DRS of NM reasoning, indicating that a higher number of rules in the context
244 often leads to more frequent mistakes.

245 **Effect on other logic datasets** Table 6 represents the accuracy comparison between LogicT5 and
246 baseline T5-large in both single-task and multi-task settings. The results indicate that training LLMs
247 on *LogicBench(Aug)* has a greater impact on logic datasets that primarily focus on logical reasoning,
248 such as FOLIO and LogicNLI. Hence, we can observe that LogicT5 consistently outperforms the
249 baseline for LogicT5 and FOLIO. However, LogiQA and ReClor encompass other forms of reasoning
250 in addition to logical reasoning, hence, LogicT5 demonstrates competitive performance on them.

251 **How do LLMs reason step-by-step?** We investigate the fraction of low-performing axioms that
252 contain various types of logical reasoning steps to predict the answer, and whether the correctness
253 of those steps is correlated with the performance. Here, we perform a case study on ChatGPT.
254 We prompt ChatGPT to generate reasoning steps along with predictions. For PL, we observe that

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

Table 6: Performance comparison between LogicT5 and baseline T5-large in terms of accuracy.

255 while the model can effectively reason the initial section of the *disjunctive syllogism* involving two
256 possibilities p or q , it encounters challenges in deducing whether q should follow from the $\neg p$. For
257 FOL, ChatGPT encounters challenges in comprehending longer logical contexts, resulting in a lack
258 of confidence in establishing the relationship between given propositions. Furthermore, to derive
259 an accurate conclusion when the rules are followed correctly, the model relies on supplementary
260 evidence. We observe that ChatGPT encounters difficulties in comprehending the nuanced meanings
261 of words such as “usually”, “normally” and “typically” when establishing sentence relationships
262 within NM reasoning. Notably, when it comes to the rule of default reasoning, ChatGPT fails to grasp
263 inherent associations between two entities that commonly share characteristics. Examples and more
264 analysis of generated explanations for each logic type are presented in Appendix E.

265 4 Related Work

266 LogiQA [11] and ReClor [19] have made notable contributions by compiling multichoice questions
267 from standardized graduate admission examinations that demand diverse forms of logical reasoning.
268 However, in contrast to our LogicBench, these datasets involve complex mixed forms of reasoning and
269 do not specifically focus on assessing logical reasoning in isolation. A few past attempts have been
270 made to create datasets to evaluate only logical reasoning while excluding other forms of reasoning.
271 For example, CLUTTER [15] covers inductive reasoning, [5] covers temporal logic, and Ruletaker
272 [4] evaluates whether a transformer-based model emulates deductive reasoning over synthetically
273 generated statements in a limited setting. LogicNLI [16] introduced a diagnostic benchmark for
274 FOL reasoning, with the dataset constructed by first automatically generating logic expressions and
275 then replacing the entity and attribute placeholders in the logic expressions with simple and random
276 subjects and predicates. FOLIO [6] gives diverse and complex logical expressions, however, it is only
277 limited to FOL. ProntoQA [14] provides explanation and reasoning steps but is limited to modus
278 ponens in FOL. Additional datasets for evaluating logical reasoning also exist such as TaxiNLI [7]
279 introduce logical taxonomy in NLI task and RuleBert [13] covers only soft logical rules. In summary,
280 LogicBench is evaluate logical reasoning in isolation and provides more diverse inference rules and
281 logic types compared to existing datasets. Extended related work is discussed in Appendix H.

282 5 Conclusions

283 To study the logical reasoning ability of LLMs, we introduced a novel benchmark called *LogicBench*
284 which consists of 25 distinct inference rules and reasoning patterns covering propositional, first-
285 order, and non-monotonic logics. We released two versions of the dataset: *LogicBench(Eval)* and
286 *LogicBench(Aug)*. *LogicBench(Eval)* serves as a high-quality, cost-effective, and reliable dataset for
287 evaluating LLMs, while *LogicBench(Aug)* can be utilized for training purposes. Through compre-
288 hensive experiments, we showed that models such as GPT-3 and ChatGPT do not perform well on
289 *LogicBench*, even though they require the application of only a single inference rule in positive (i.e.,
290 label ‘Yes’) data instance. Furthermore, we demonstrated that LLMs trained using *LogicBench(Aug)*
291 showcase an improved understanding of logical reasoning, resulting in a better performance on
292 existing logic datasets. Though *LogicBench* facilitates the evaluation and improvement of the logical
293 reasoning ability of LLMs, it can be further extended by incorporating other inference rules and logic
294 types; and having data instances that require applications of multiple inference rules.

295 References

- 296 [1] Sajjad Beygi, Maryam Fazel-Zarandi, Alessandra Cervone, Prakash Krishnan, and Siddhartha
297 Jonnalagadda. Logical reasoning for task oriented dialogue systems. In *Proceedings of the*
298 *Fifth Workshop on e-Commerce and NLP (ECNLP 5)*, pages 68–79, Dublin, Ireland, May 2022.
299 Association for Computational Linguistics.
- 300 [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
301 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
302 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
303 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
304 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
305 Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle,
306 M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information*
307 *Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- 308 [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
309 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
310 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 311 [4] Peter Clark, Oyvind Tafjord, and Kyle Richardson. Transformers as soft reasoners over language.
312 In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences*
313 *on Artificial Intelligence*, pages 3882–3890, 2021.
- 314 [5] Christopher Hahn, Frederik Schmitt, Jens U. Kreber, Markus Norman Rabe, and Bernd
315 Finkbeiner. Teaching temporal logics to neural networks. In *International Conference on*
316 *Learning Representations*, 2021.
- 317 [6] Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy
318 Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, et al. Folio: Natural language reasoning
319 with first-order logic. *arXiv preprint arXiv:2209.00840*, 2022.
- 320 [7] Pratik Joshi, Somak Aditya, Aalok Sathe, and Monojit Choudhury. TaxiNLI: Taking a ride
321 up the NLU hill. In *Proceedings of the 24th Conference on Computational Natural Language*
322 *Learning*, pages 41–55, Online, November 2020. Association for Computational Linguistics.
- 323 [8] Daniel Khashabi. *Reasoning-Driven Question-Answering for Natural Language Understanding*.
324 University of Pennsylvania, 2019.
- 325 [9] Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark,
326 and Hannaneh Hajishirzi. UNIFIEDQA: Crossing format boundaries with a single QA system.
327 In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1896–1907,
328 Online, November 2020. Association for Computational Linguistics.
- 329 [10] Vladimir Lifschitz. Benchmark problems for formal nonmonotonic reasoning: Version 2.00.
330 In *Non-Monotonic Reasoning: 2nd International Workshop Grassau, FRG, June 13–15, 1988*
331 *Proceedings 2*, pages 202–219. Springer, 1989.
- 332 [11] Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: a
333 challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings*
334 *of the Twenty-Ninth International Conference on International Joint Conferences on Artificial*
335 *Intelligence*, pages 3622–3628, 2021.
- 336 [12] OpenAI. Gpt-4 technical report, 2023.
- 337 [13] Mohammed Saeed, Naser Ahmadi, Preslav Nakov, and Paolo Papotti. RuleBERT: Teaching
338 soft rules to pre-trained language models. In *Proceedings of the 2021 Conference on Empirical*
339 *Methods in Natural Language Processing*, pages 1460–1476, Online and Punta Cana, Dominican
340 Republic, November 2021. Association for Computational Linguistics.

- 341 [14] Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal anal-
 342 ysis of chain-of-thought. In *The Eleventh International Conference on Learning Representations*,
 343 2023.
- 344 [15] Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. CLUTRR:
 345 A diagnostic benchmark for inductive reasoning from text. In *Proceedings of the 2019 Confer-*
 346 *ence on Empirical Methods in Natural Language Processing and the 9th International Joint*
 347 *Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4506–4515, Hong Kong,
 348 China, November 2019. Association for Computational Linguistics.
- 349 [16] Jidong Tian, Yitian Li, Wenqing Chen, Liqiang Xiao, Hao He, and Yaohui Jin. Diagnosing the
 350 first-order logical reasoning ability through LogicNLI. In *Proceedings of the 2021 Conference*
 351 *on Empirical Methods in Natural Language Processing*, pages 3738–3747, Online and Punta
 352 Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- 353 [17] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,
 354 Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan
 355 Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson,
 356 Kirby Kuznia, Krима Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir
 357 Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh
 358 Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A,
 359 Sumanta Patro, Tanay Dixit, and Xudong Shen. Super-NaturalInstructions: Generalization
 360 via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on*
 361 *Empirical Methods in Natural Language Processing*, pages 5085–5109, Abu Dhabi, United
 362 Arab Emirates, December 2022. Association for Computational Linguistics.
- 363 [18] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
 364 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *ICLR*,
 365 2021.
- 366 [19] Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension
 367 dataset requiring logical reasoning. In *International Conference on Learning Representations*.
- 368 [20] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
 369 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*
 370 *preprint arXiv:2303.18223*, 2023.

371 **Paper Checklist**

372 **For all authors...**

- 373 1. Do the main claims made in the abstract and introduction accurately reflect the paper’s
 374 contributions and scope?
 375 Yes
- 376 2. Have you read the ethics review guidelines and ensured that your paper conforms to them?
 377 Yes
- 378 3. Did you discuss any potential negative societal impacts of your work?
 379 No, we do not expect negative societal impacts as a direct result of the contributions in our
 380 paper
- 381 4. Did you describe the limitations of your work?
 382 Yes, refer to Section 5.

383 **If you are including theoretical results...**

384 1. Did you state the full set of assumptions of all theoretical results?

385 N/A

386 2. Did you include complete proofs of all theoretical results?

387 N/A

388 **If you ran experiments...**

389 1. Did you include the code, data, and instructions needed to reproduce the main experimental
390 results (either in the supplemental material or as a URL)?

391 Yes, the anonymous URL is at the end of the abstract.

392 2. Did you specify all the training details (e.g., data splits, hyperparameters, how they were
393 chosen)?

394 Yes, refer to Section 3.1 and Appendix D.

395 3. Did you report error bars (e.g., with respect to the random seed after running experiments
396 multiple times)?

397 Yes, we reported the average results across three prompts (refer to Section 3.1).

398 4. Did you include the amount of compute and the type of resources used (e.g., type of GPUs,
399 internal cluster, or cloud provider)?

400 Yes, refer to Appendix D.

401 **If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...**

402 1. If your work uses existing assets, did you cite the creators?

403 Yes, refer to Section 1, Section 3, and Appendix D.

404 2. Did you mention the license of the assets?

405 Yes, refer to Appendix D.

406 3. Did you include any new assets either in the supplemental material or as a URL?

407 Yes

408 4. Did you discuss whether and how consent was obtained from people whose data you're
409 using/curating?

410 N/A

411 5. Did you discuss whether the data you are using/curating contains personally identifiable
412 information or offensive content?

413 Yes, the collected data does not contain personally identifiable information or offensive
414 content.

415 **If you used crowdsourcing or conducted research with human subjects...**

416 1. Did you include the full text of instructions given to participants and screenshots, if applica-
417 ble?

418 N/A

419 2. Did you describe any potential participant risks, with links to Institutional Review Board
420 (IRB) approvals, if applicable?

421 N/A

422 3. Did you include the estimated hourly wage paid to participants and the total amount spent
423 on participant compensation?

424 N/A