Logical forms complement probability in understanding language model (and human) performance

Anonymous ACL submission

Abstract

With the increasing interest in using large language models (LLMs) for planning in natural language, understanding their behaviors becomes an important research question. This work conducts a systematic investigation of LLMs' ability to perform logical reasoning in natural language. We introduce a controlled dataset of hypothetical and disjunctive syllogisms in propositional and modal logic and use it as the testbed for understanding LLM performance. Our results lead to novel insights in predicting LLM behaviors: in addition to the probability of input (Gonen et al., 2023; McCoy et al., 2024), logical forms should be considered as important factors. In addition, we show similarities and discrepancies between the logical reasoning performances of humans and LLMs by collecting and comparing behavioral data from both.

1 Introduction

001

006

007 008

011

012

019

024

027

Logical reasoning is a fundamental aspect of building AI systems for reliable decision-making (Kautz et al., 1992, inter alia)-given a set of premises, an AI system should be able to deduce valid conclusions. With the advent of large language models (LLMs; Touvron et al., 2023; Jiang et al., 2023; AI@Meta, 2024, inter alia), there has been a surge of interest in using these models to assist planning and decision-making (Huang et al., 2022, inter alia); therefore, understanding the logical reasoning capabilities becomes crucial in understanding the reliability and potential of LLMs in planning. While recent work has shown that LLMs exhibit decent performance on logical reasoning problems (Liu et al., 2020; Ontanon et al., 2022; Wan et al., 2024, inter alia), there is still a lack of fine-grained understanding of the logical formsamong many argument forms presented in natural language (Shieber, 1993), do LLMs perform equally well, or do they exhibit preferences for



Figure 1: Illustration of the fact that perplexity does not serve as a reliable indicator of logical reasoning performance; and therefore, neither does probability. The distributions of the probabilities assigned to the ground-truth answer (i.e., soft accuracy; Y-axis) by Llama-3-70B are plotted against the perplexity of the corresponding example question (X-axis) and grouped by (a) modality, (b) argument forms, and (c) logic interpretation content. Each group consists of 20 randomly selected examples with other factors controlled.

certain argument forms? Do more complex components of logical forms, such as modalities, matter for LLM performance?

In this work, we investigate the logical reasoning capabilities of LLMs by assessing their performance on different logical forms. We curate a dataset of natural language statements and questions based on several logical forms in both proposi-

047

tional and modal logic, which is designed to mirror 049 reasoning in daily communication. An example is shown in §3.3. We then conduct a series of controlled experiments to analyze the performance of a set of LLMs on the dataset. Although our findings generally align with those by Gonen et al. (2023) and McCoy et al. (2024), who suggest that LLMs excel on examples with high probability, our results indicate that logical form, including but not limited 057 to modalities and argument forms, is a crucial complementary factor in predicting the performance of LLMs (Figure 1). Additionally, with meaningful real-world interpretations, we find that: 061

- 1. LLMs are still far from being perfect in atomiclevel propositional and modal logic reasoning.
 - 2. LLMs prefer an affirmative answer under the modality of possibility, whereas they prefer a negative answer under the modality of necessity.
 - 3. In line with the recent results on categorical syllogisms (Eisape et al., 2024), we verify on hypothetical and disjunctive syllogisms that LLMs achieve better performance on certain logical forms that humans perform well. However, some logical forms receive favor from LLMs, while the phenomena lack support from human intuition or human behavioral data.

This paper is structured as follows. After reviewing related work (§2), we describe the dataset synthesis process (§3). We report the LLM reasoning results on our data (§4), compare them with human performance (§5). We conclude by discussing the implications of our results and the limitations(§6).

2 Related Work

065

081

084

094

095

Logical reasoning benchmarks. Prior LLM logical reasoning benchmarks (Liu et al., 2020; Han et al., 2022, inter alia) focus on complex, multihop reasoning problems with manually annotated problems, making cross-problem comparisons challenging. Recent work has introduced benchmarks with synthesized natural-language questions using predefined logical formulas and substitution rules (Saparov and He, 2022; Saparov et al., 2023; Parmar et al., 2024; Wan et al., 2024, inter alia). Compared to them, our work uniquely incorporates modal logic, which has been largely unexplored in existing benchmarks-while Holliday and Mandelkern (2024) present a case study, our approach offers two key advances: controlled knowledge bias in logic interpretations (§3.3) and a more rigorous statistical evaluation framework (§4.1).

Propositional and modal logic reasoning in language models. Recent work has explored training and finetuning language models specifically for logical reasoning (Clark et al., 2021; Hahn et al., 2021; Tafjord et al., 2022). Our work differs in two key aspects: (1) we evaluate general-purpose language models through prompting, a cost-efficient setup that has been widely adopted in recent years, and we focus on propositional and alethic modal logic rather than temporal (Hahn et al., 2021) or epistemic (Sileo and Lernould, 2023) logic; ¹ (2) unlike studies comparing LLM and human performance on categorical syllogisms (Eisape et al., 2024, *inter alia*),² we focus on hypothetical and disjunctive syllogisms with considerations of modality. 099

100

101

102

103

104

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

Human logic reasoning. Work on human reasoning capabilities has informed studies of LLM logical reasoning: Eisape et al. (2024) compared LLM syllogistic reasoning with human behavior results (Ragni et al., 2019) under the framework of the Mental Models Theory (Johnson-Laird, 1983); Lampinen et al. (2024) found similar content effects in human and LLM reasoning, supporting the need to control for common-sense knowledge in benchmarks (§3.2); Belem et al. (2024) studied human and LLM perception of uncertainty at a lexical level. Compared to them, we focus on the propositional and modal logic reasoning process and contribute new behavioral data.

3 Dataset

We curate a dataset of natural-language multichoice questions to measure the logical inference performance of LLMs. Starting from propositional and modal logical forms as templates (§3.1), we assign meanings (e.g., real-world interpretations) to each variable and translate templates into naturallanguage Yes/No questions ((§3.2). A subsidary visualization of the process is shown in Figure A1.

3.1 Background: Propositional and Modal Logic

Propositional logic studies the relation between propositions. In this framework, each proposition is typically represented by a variable, and multi-

¹Technically, any logic that involves non-truth-functional operators, including first-order logic, temporal logic, and epistemic logic, can be viewed as a modal logic; however, we adopt the most restrictive sense of *modal logic* (Ballarin, 2023) and use it interchangeably with *alethic modal logic*.

²We refer readers to Zong and Lin (2024) for a more comprehensive review of categorical syllogisms.

ple propositions combine with logical connectives (e.g., \lor and \rightarrow) to form compound propositions.

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

170

172

173

174

175

177

178

179

180

182

184

185

187

In propositional logic, a proposition can be evaluated as either true or false; however, this system can be overly simplistic when dealing with the complexity of real-world events. Consider the statement Alice is not eating, while it is true in a world where Alice is not eating, it may become false in a hypothetical possible world where Alice is indeed eating. This idea, known as possible world semantics (Kripke, 1959), provides a framework for more nuanced statements about event possibilities, such as *Alice may be eating* and *Alice must be* eating. The former statement can be understood as there exists a possible world where Alice is eating, and the latter can be understood as in all possible worlds, Alice is eating.³ Normal modal logic (Kripke, 1963) formalizes this idea and extends propositional logic to reason about event necessity and possibility. In the Backus-Naur form, a normal modal logic system \mathcal{L} can be written as

$$\mathcal{L}: \varphi \coloneqq p \mid \neg \varphi \mid \Box \varphi \mid \Diamond \varphi \mid \varphi \lor \varphi \mid \varphi \land \varphi \mid \varphi \to \varphi, \qquad (1)$$

where p is a propositional variable that serves as an atom in \mathcal{L} , \neg is the negation operator, \Box is the necessity operator (*must*), \diamond is the possibility operator (*may*), \lor is logical disjunction (*or*), \land is logical conjunction (*and*), and \rightarrow is the logical implication operator (*if...then*). φ denotes the syntactic category of a formula in \mathcal{L} . The right-hand side of Eq. (1) describes all possible logical formulas under the system \mathcal{L} : for example, if $\varphi \in \mathcal{L}$, the rules imply that $\neg \varphi \in \mathcal{L}$, $\Box \varphi \in \mathcal{L}$, and so on. Following the convention in logic, the operator precedence is $\{\neg, \Box, \diamondsuit\} \succ \{\lor, \land\} \succ \{\rightarrow\}$.

Indeed, the operators $(\neg, \Box, \rightarrow)$ forms a functional complete set of operators under \mathcal{L} . Suppose φ and ψ are variables that represent logical formulas. The logical or (\lor) and logical and (\land) operators can be rewritten with logical not (\neg) and logical implication (\rightarrow) , as follows:

$$\begin{aligned} \varphi \lor \psi \Leftrightarrow \neg \varphi \to \psi, \\ \varphi \land \psi \Leftrightarrow \neg (\varphi \to \neg \psi). \end{aligned}$$
(2)

Possibility operator \diamond can also be derived from the necessity operator.

$$\rangle \varphi \Leftrightarrow \neg \Box \neg \varphi \tag{3}$$

Deduction and sequent. Given a formula set Γ as premises, if a deduction to a conclusion φ exists using axiom schemata and inference rules under the normal modal logic, we say the premises *infer* the conclusion, and the deduction can be represented as a logic *sequent* $\Gamma \vdash \varphi$. If a formula set Γ do not infer the conclusion, we denote it as $\Gamma \nvDash \varphi$ and call it a *non-entailment*. 188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

230

231

232

233

234

235

3.2 Translating Logic to Natural Language

An *interpretation* maps propositional variables to concrete meanings. For example, under the interpretation that p is "Jane is eating apples" and q is "John is eating oranges", the logical formula $p \lor q$ becomes "Jane is eating apples or John is eating oranges."

Choices of interpretation, i.e., the concrete content of the sentence, should not affect the underlying logical reasoning process. However, in natural-language utterances, reasoning can be influenced by various confounding factors. Knowledge bias is a common pitfall. For example, given the logical form $\{\Box p \rightarrow \Box \neg q, \Box p\} \vdash \Box \neg q$, regardless of p's interpretation, if we interpret $\neg q :=$ "*Cats are not animals*" then the conclusion will be "It is certain that cats are not animals." But common-sense knowledge suggests that "it is *certain that cats are animals*" ($\Box q$), which logically contradicts the existing premise set.⁴ Such bias will complicate logical reasoning (Lampinen et al., 2024) and should be avoided in data curation. Besides, each variable should have independent interpretation, as detailed in Appendix B.2.

After being assigned interpretations, each logical form is further articulated as a yes-no question on whether the conclusion can be inferred from the premises. To mitigate the ambiguity in natural language, we design heuristic rules to translate logic forms into less ambiguous English, which are detailed in Appendix B.2. For the exact wordings we used, see Table A1 in Appendices. If a valid deduction exists (\vdash) for the logical form, the ground truth answer is Yes, otherwise No. The answer is solely determined by the logical form and is independent of the interpretation.

3.3 Involved Logical Forms

Translated logical forms can have varying degrees of naturalness. For example, the *necessitation rule* $\{\varphi\} \vdash \Box\varphi$, which translates to " φ is true; therefore,

³The possible world semantics, therefore, connects the notion of *necessity* and *possibility* to the universal and existential quantification (\forall, \exists) under first-order logic.

⁴This confounding factor affects the examples in Table 18 of Han et al. (2022).

274

275

291

292

293

294

295

296

297

388

301

303

304

288

236

237

241

242

243

245

246

247

248

249

251

252

263

266

269

270

271

272

273

$$\{\varphi \lor \psi, \neg \varphi\} \vdash \psi, \qquad (\lor^{\mathrm{L}})$$
$$\{\neg \varphi \lor \psi, \neg \varphi\} \vdash \psi \qquad (\to^{\mathrm{L}}: \text{ modus ponens})$$

$$\begin{cases} \neg \varphi \to \psi, \neg \varphi \} \vdash \psi, \quad (\to^{\mathrm{L}}; \text{ modus ponens}) \\ \{ \varphi \lor \psi, \neg \psi \} \vdash \varphi, \qquad (\lor^{\mathrm{R}}) \end{cases}$$

 $\{\neg \varphi \rightarrow \psi, \neg \psi\} \vdash \varphi.$ (\rightarrow^{R} ; modus tollens)

Despite the semantic similarity, these logical forms translate to different natural-language questions. For example, taking the interpretations of φ := Jane is watching a show and ψ := *John is reading a book*, \vee^{L} translates to

Consider the following statements:

Jane is watching a show or John is reading a book.

Jane isn't watching a show.

Ouestion: Based on these statements, can we infer that John is reading a book?

With the same interpretation, \rightarrow^{L} 's translation of the first statement is If Jane isn't watching a show, then John is reading a book.

According to the commutativity of disjunction operator, we group \vee^{L} and \vee^{R} together as *disjunc*tive syllogism, alongside two hypothetical syllogism groups, modus ponens (\rightarrow^{L}) and modus tollens (\rightarrow^{R}). All the logical forms shown above are valid sequents with ground-truth answer Yes. To balance the dataset, we introduce some logic fallacies that generate questions with ground-truth label No. By flipping the second premises and the conclusions, we obtain the following fallacies:

$$\{\varphi \lor \psi, \psi\} \nvDash \neg \varphi, \qquad (\lor_{\Bbbk}^{\mathsf{I}})$$

$$\{\varphi \lor \psi, \varphi\} \nvDash \neg \psi,$$

$$\{\neg\varphi \to \psi, \varphi\} \nvDash \neg\psi, \qquad (\to^{\mathrm{R}}_{\varkappa})$$

where $\vee^{\mathrm{L}}_{\nvDash}$ and $\vee^{\mathrm{R}}_{\nvDash}$ are grouped as *affirming the dis-junction*, $\rightarrow^{\mathrm{L}}_{\nvDash}$ and $\rightarrow^{\mathrm{R}}_{\nvDash}$ corresponds to *affirming* the consequent and denying the antecedent, respectively. In our dataset, we require the formulas φ and ψ to the form of $\mathfrak{M}p$ and $\mathfrak{M}q$, where p and q are propositional variables, each assigned with an interpretation. Both variables are constrained under the same modality \mathfrak{M} , which can be necessity (\Box) ,

possibility (\diamond) or no modality (\varnothing). Pairing with four rules and theorem-fallacy variations, we have a total of $3 \times 4 \times 2 = 24$ forms.

3.4 Involved Logic Interpretations

For logic interpretations, we generate a set of verb phrases by prompting the CodeLlama 2 model (Rozière et al., 2024), and select 204 of them manually. and combine them with top-200 popular baby names in the US into subject-verb-object pairs,⁶ such as (*Ray, make, a pizza*). We randomly generate 1000 interpretations with two pairs each. The same set of interpretations is applied to variables p, q in each logic sequent's natural langauge template. In total, there are $24 \times 1000 = 24000$ question, with samples shown in Table A1.

Experiment 4

4.1 **Metrics and Investigated Models**

Hu and Levy (2023) have suggested that the standard approach of greedily decoding yes-no strings (Dentella et al., 2023) may underestimate the competence of a language model; therefore, we adopt a probability-based metric to evaluate the model performance. In our evaluation protocol, the predicted likelihood of the tokens Yes and No, conditioned on the prompt s—denoted as p(Yes | s) and p(No | s), respectively-serve as the soft labels for yes-no answers. The soft accuracy \hat{p} on the single example with ground-truth answer $y \in \{\text{Yes}, \text{No}\}$ is defined as the relative probability of y:

$$\hat{p} = \frac{p(\mathsf{No} \mid s)\mathbb{1}[y = \mathsf{No}] + p(\mathsf{Yes} \mid s)\mathbb{1}[y = \mathsf{Yes}]}{p(\mathsf{No} \mid s) + p(\mathsf{Yes} \mid s)},$$

where $\mathbb{1}[\cdot]$ is the indicator function that returns 1 if the condition is true and 0 otherwise. This relative probability can also be viewed as the confidence score of the model on the ground-truth answer. The soft accuracy Acc_{soft} of a model on the entire dataset \mathcal{D} is defined as the average soft accuracy over all examples,

$$Acc_{soft} = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \hat{p}_i.$$
 29

We use a zero-shot setting to investigate the general performance of the models' logical inference capabilities—while adding detailed instructions or few-shot demonstrations may increase the absolute performance, they are at the cost of introducing

⁵Nevertheless, we report the experiment results on necessitation rule in Appendix C.1.

⁶https://www.ssa.gov/oact/babynames/names.zip

	Ove	rall	Leader	board		Modality	/			Argume	nt Form		
Model	(Ra	nk)	(Ra	nk)	Ø		\diamond	$\nabla^{\mathrm{L,R}}_{\vdash}$	$\rightarrow^{\mathrm{L}}_{\vdash}$	$\rightarrow^{\rm R}_{\vdash}$	$\vee^{\mathrm{L,R}}_{\nvDash}$	$\rightarrow^{\mathrm{L}}_{\not\vdash}$	$\rightarrow^{\mathrm{R}}_{\nvDash}$
mistral-7b mistral-8x7b	0.645	(4) (1)	0.145	(7) (5)	0.464	0.496	0.974 0.874	0.877 0.963	0.663	0.280	0.434	0.653	0.939 0.813
llama-2-7b	0.335	(10)	0.094	(10)	0.262	0.207	0.538	0.444	0.147	0.315	0.208	0.451	0.468
llama-2-13b	0.513	(9)	0.110	(9)	0.488	0.362	0.688	0.418	0.581	0.393	0.631	0.436	0.591
llama-2-70b	0.611	(5)	0.127	(8)	0.616	0.471	0.746	0.446	0.845	0.518	0.775	0.389	0.694
llama-3-8b	0.565	(6)	0.239	(3)	0.598	0.460	0.639	0.526	0.470	0.332	0.664	0.625	0.716
llama-3-70b	0.714	(2)	0.362	(1)	0.745	0.554	0.843	0.606	0.773	0.515	0.882	0.661	0.788
yi-34b	0.518	(8)	0.226	(4)	0.457	0.413	0.683	0.346	0.498	0.205	0.685	0.638	0.737
phi-2	0.532	(7)	0.155	(6)	0.469	0.456	0.673	0.670	0.757	0.522	0.365	0.402	0.510
phi-3-mini	0.690	(3)	0.272	(2)	0.657	0.536	0.877	0.839	0.974	0.475	0.664	0.462	0.604
OpenAI-o1	0.926	N/A	N/A	N/A	1.000	0.773	1.000	0.895	1.000	0.775	0.919	1.000	1.000
Gemini-1.5-Pro	0.859	N/A	N/A	N/A	0.831	0.748	0.997	1.000	1.000	0.919	0.661	0.991	0.638
human	0.595	N/A	N/A	N/A	0.589	0.566	0.640	0.691	0.901	0.628	0.594	0.225	0.411

Table 1: Overall and break-down accuracies of different models, as well as their HuggingFace OpenLLM Leaderboard performance and relative ranking (Fourrier et al., 2024). Each argument form category denotes the union of the fine-grained categories specified in the superscripts and subscripts—for example, $\bigvee_{\vdash}^{L,R}$ denotes the entire disjunctive syllogism group. **Boldfaced** values indicate the row-wise maximum for each factor. Note that due to technical limitations of commercial LLMs, results from OpenAI-o1 (OpenAI, 2024) and Gemini-1.5-pro (Team et al., 2024) are greedy-decoding based evaluation on 2,000 random samples that serve as references, and are therefore not directly comparable to other probability-based evaluations. Human results are detailed in §5.

possibly undesired confounding factors or behaviors, such as simply copy-pasting the answers in the examples.

306

310

311

312

313

314

315

316

317

319

320

321

323

325

326

327

We evaluate on the following models with opensourced weights: mistral-7b-v0.2 and -8x7b (Jiang et al., 2023, 2024); llama-2-7b, -13b and -70b (Touvron et al., 2023); 3.1 version of llama-3-8b and -70b (AI@Meta, 2024); yi-34b (01.AI, 2024); phi-2 and phi-3-mini (Microsoft, 2023, 2024).⁷

4.2 Results: Performance w.r.t. Logical Forms

We evaluate the aforementioned models with the probability-based protocol (Table 1). Generally, models that rank higher in the leaderboard also achieve higher soft accuracy on our dataset. The break-down accuracies on modalities and argument forms reveal that:

- 1. (Modality) All models consistently perform better on the possibility (◊) than necessity (□) or plain propositional logic.
- (Argument Forms) The pattern is more diverse, yet most of the models struggle the most on modus tollens (→^R_⊢) within logic sequents (i.e., questions with ground-truth answers Yes), and affirming the consequent (→^L_⊬) within fallacies.

4.2.1 Analysis on Logic Sequents

To systematically analyze the effect on model performance of each factor of interest, as well as crossvalidating the observations above, we fit a linear mixed-effects model (Raudenbush, 2002) to the soft accuracy data on valid logic sequents (i.e. with ground truth of Yes) across different LLMs and logical forms, 329

331

332

333

334

335

336

337

339

340

341

342

343

344

345

346

348

349

350

351

352

353

$$Acc_{soft} \sim Modality + ArgForm + Perplexity$$

+ (1 + Perplexity | LLM), (4)

with the linear fixed effects of (i.) modality, (ii.) argument form, and (iii.) input perplexity. Individual probability, coupled with a constant term, is modeled as a random effect to account for potential model-specific biases. Here, *Perplexity* denotes the perplexity of the input text $(x_1x_2...x_N)$, which is defined as the exponential of the token-wise average negative log-likelihood of the text given a specific language model:

Perplexity = exp
$$\left(-\frac{1}{N}\sum_{i=1}^{N}\log p(x_i \mid x_{< i})\right)$$

The mixed-effects model yields a marginal R^2 of 0.342 and a conditional R^2 of 0.543, suggesting a reasonable predictive power. The likelihood ratio test on the full regression model vs. the null regression model without each of the fixed effects

⁷Our evaluation protocol technically requires the conditional probabilities of specified answers given a prompt, which are not supported by most commercial models; however, we report the greedy-decoding accuracy of these models on a sample subset for reference.



Figure 2: Estimated marginal means of logical form factors in the mixed-effects model of Eq. (4), along with their 95% confidence intervals.

Hypothesis	p-value
propositional < may must < propositional must < may	$< 0.001 \\ < 0.001 \\ < 0.001$
disjunctive < modus ponens modus tollens < modus ponens modus tollens < disjunctive	< 0.001 < 0.001 < 0.001

Table 2: Hypothesis testing results on the effect of logical form factors on soft accuracy (Figure 2).

yields a significant result (p < 0.001), suggesting the importance of all these factors in determining the model performance.

354

355

357

359

367

371

Fixed effects. In line with Gonen et al. (2023) and McCoy et al. (2024), we find a negative correlation between perplexity and soft accuracy (p < 0.001); however, the correlation is weak ($\rho = -0.09$), which suggests the necessity of the complementary factors below in predicting LLM performance.

For different modalities and argument forms, we estimate their marginal means on soft accuracy (Figure 2), and perform pairwise hypothesis testing on the estimated coefficients (Table 2). The results generally align with the general observations on the full dataset. The only exception is that modus ponens (\rightarrow^{L}), instead of disjunctive syllogism (\lor), appears to be the easiest argument form (i.e., the one with the highest soft accuracy) among all.

Random effects. We analyze the per-LLM random effects on the soft accuracy (Figure 3). All
the model-specific mixed effects of perplexity are
negative, suggesting the negative correlation be-



Figure 3: Illustration of per-model random effects on soft accuracy in the mixed-effects model of Eq. (4) with 99.9% confidence intervals. (a) Mixed effects (i.e., the sum of fixed and random effects) of perplexity. (b) Intercept random effects (i.e., constant term per model on soft accuracy), with the model performance rank (Table 1) annotated in parentheses.

tween perplexity and soft accuracy is consistent across models (Figure 3a). While the intercept random effects are not perfectly aligned with the model performance—since the perplexity random effects may introduce confounding factors—higherranked models generally tend to have higher intercept random effects (Figure 3b), which crossvalidates the general performance ranking. 376

377

378

379

380

381

382

384

385

386

387

388

390

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

4.2.2 Extended Analysis on the Negative Perplexity–Performance Correlation

We further investigate the negative correlation between perplexity and model performance through a controlled experiment: we create a mirror dataset of the same size, keeping all the logical formulas while interpreting them with nonsensical words. For example, the formula $\Diamond(\varphi \lor \psi)$ may be interpreted as *it's possible that Neva is <u>balaring</u> a <u>montery or Lucille is sweeling prandates</u>, where the underlined words and phrases are nonsensical. Intuitively, the perplexity of the problems in this mirror dataset should be much higher than that of the primary dataset problems (§3) under any reasonably trained language model.*

We analyze the correlation between perplexity and model performance (Figure 4). As desired, the perplexity of problems with nonsensical words are indeed much higher than that of the primary dataset (≈ 30 vs. ≈ 10). The significant portion of horizontal and inclined lines in the figures again suggests that perplexity is not a reliable predictor of model performance. Meanwhile, the overall



Figure 4: Correlation between mean perplexity and mean confidence score on each logic sequent. Each point represents an average over a group of 1000 prompts that share the same underlying logic sequent. Two connected dots share the same logic formula.



Figure 5: Estimated marginal means of the factors in the mixed-effects model of Eq. (5) with 95% confidence intervals. Higher coefficients indicate a higher tendency to affirm the claim.

parallelism of the lines echos our results that logical forms are important factors for such prediction.

407

408

409

410

411

412

413

414

415

416

417

418

419

4.2.3 The Affirmation Bias over Modalities

One key argument of Dentella et al. (2023) is that large language models exhibit a bias towards affirming the claim, i.e., answering Yes more frequently than No. We investigate this phenomenon by fitting a mixed-effects model

$$\frac{P(\operatorname{Yes} \mid s)}{P(\operatorname{Yes} \mid s) + P(\operatorname{No} \mid s)} \sim Modality + ArgForm + Perplexity + (1 + Perplexity \mid LLM),$$
(5)

which has the same structure as Eq. (4), except the dependent variable being the relative probability of answering Yes conditioned on input text s.

We present the estimated marginal means of
the factors in the mixed-effects model (Figure 5).
While our results confirm the affirmation bias on
propositional logic, such bias is slightly less pronounced on the possibility modality (◊, around
0.03), and the models even show a bias towards
rejecting claims under the necessity modality (□).

5 Human Experiments

LLMs are trained on text produced by humans and are able to generate plausible text; therefore, there have been interests in using LLMs as human models (Eisape et al., 2024; Misra and Kim, 2024, *inter alia*). Following this line of work, we conduct a human behavioral experiment to ground the LLM reasoning behavior. Using samples from our primary dataset, we collected 710 responses from adults fluent in English through Prolific.⁸ More experiment details can be found in Appendix A.2. 427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

The average human accuracy on each group is shown in the last row of Table 1.⁹ Aligned with our LLM results (§4), on modalities the overall human results also show an accuracy order of ($\diamond \succ \varnothing \succ \Box$), and on argument forms, modus ponens (\rightarrow^{L}) is the most accurately answered pattern.

To further investigate the interactions of logic factors, we fit a generalized linear mixed-effects model (Bates et al., 2015) to verify the effect of modality and argument forms on human logic reasoning accuracy (Eq. (6) and Figure 6).

$$logit(Acc) \sim Modality + ArgForm + Rt + (1 + Rt | ParticipantID), \quad (6)$$

where Acc is the binary accuracy of human responses, and Rt is the response time. The generalized mixed-effects model yields a marginal R^2 of 0.121 yet a 0.419 conditional R^2 , indicating a diverse response pattern across participants. The likelihood ratio test on the full model against the null model shows that only the effect of argument form

⁸https://prolific.com

⁹Human responses are binary classes, so correct and incorrect responses are coded as 1 and 0, respectively.



Figure 6: Estimated marginal means of logical form factors in the generalized mixed-effects model of Eq. (6), along with their 95% confidence intervals.

is significant ($\chi^2(2) = 25.6$, p < 0.001). However, in accordance with the overall performance, we find modus ponens (\rightarrow^L) has a significantly higher effect than other two valid argument forms. This confirms that logical forms can also have a significant impact on human reasoning accuracy, which is consistent with the LLM results, although the effect sizes are not the same.

6 Conclusion and Discussion

We present an analysis of hypothetical and disjunctive syllogisms on propositional and modal logic and systematically analyze the LLM performance on the dataset. Our analysis provides novel insights on explaining and predicting LLM performance: in addition to the perplexity or probability of the input text, the underlying logic forms play an important role in determining the performance of LLMs. In addition, we compare the behaviors of LLMs and humans using the same data through human behavioral experiments. We discuss the implications of our results as follows.

Probability in language models. Probability and perplexity are often used as intrinsic evaluation metrics for language models. While Gonen et al. (2023) and McCoy et al. (2024) show that probability and perplexity correlate well with LLM performance, literature in program synthesis with LLMs shows little correlation between probability and execution-based evaluation results (Li et al., 2022; Shi et al., 2022). This work does not necessarily contradict either line but rather provides complementary factors for analyzing LLM performance.

We argue that probability may have become an overloaded term in analyzing LLMs. Low probability may be due to one or more of the following non-exhaustive reasons: (1) out-of-context content, (2) ungrammatical language, or (3) grammatical but semantically awkward content (cf. the mirror dataset in §4.2.2), (4) reasonable but rare content. We hypothesize that the probability of language models may not be essentially able to capture all these nuanced differences, and call for encoding and decoding algorithms—such as Meister et al. (2023)—that can better decompose the probability into finer-grained and explainable components. 490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

Comparing humans and LLMs. What is our goal for building LLMs? To achieve better performance on practical tasks or to build a more human-like model? Our results, together with Eisape et al. (2024), suggest that these two goals may not be perfectly aligned by revealing a mixture of similarity and discrepancy between LLMs and humans-for example, while LLMs exhibit higher benchmark performance than humans on our dataset and show the same argument form preferences with humans (Figures 2 and 6), they also show systematic biases that we do not find significant in human reasoning (e.g., disfavoring the necessity modality, §4.2.3). While there has been positive evidence of using LLMs as human models in psycholinguistic studies (Misra and Kim, 2024, inter alia), our results suggest executing such approaches cautiously.

On the relation between modality and performance. Our results show that there is a significant difference in performance between necessity and possibility modalities, with the former much lower than the latter (Table 1). Part of the reason for this is that LLMs have a significant tendency to say "No" to necessity modality (Figure 5).

On the one hand, our results extend the conclusion of Dentella et al. (2023) that LLMs generally respond positively—LLM behaviors may be significantly affected by finer-grained factors, including but not necessarily limited to the modality involved in the input. On the other hand, while LLMs systematically tend to answer "No" to questions in necessity modality, we do not find related evidence in human experiments, which leads us to hypothesize that such rejection bias comes from either the model architecture or the training strategies, such as the reinforcement learning with human feedback (RLHF; Ouyang et al., 2022) protocol. We leave this as an open question for future research.

458

459

Limitations

541

542

555

561

562

567

571

574

579

580

582

583

584

585

This work comes with two major limitations:

1. While we have verified that our data has a low 543 perplexity $(9.82 \pm 2.47 \text{ under mistral-7b; much})$ 544 lower than that of the data by Wan et al. (2024), 545 25.44), and, therefore, are similar enough to natural language utterances, the synthetic language 547 cannot fully substitute natural language in daily 548 life. Our dataset and analysis are not comprehensive enough to cover many nuanced examples that may appear in real communication, especially when context-dependent understanding is crucial to conveying communication goals. 553

> 2. Despite more than 7,000 languages worldwide, as a first step, our material only covers English. This narrow focus is due to the languages the authors are proficient in and the coverage of the language models. We acknowledge the importance of extending the scope of this work to a more comprehensive set of languages and leave the extension as an immediate follow-up step.

In addition, the sample size of human experiments is somewhat limited. We leave more comprehensive human behavioral data collection and analysis to future work.

Ethics Statement

While this work involves human logical reasoning experiments, we have ensured that (1) the data are generated procedurally following templates listed in the paper and (2) there is no harmful content in the atomic logical interpretations, reviewed by all the authors. In addition, we have ensured that all participants are paid a fair wage through the Prolific platform. Instructions and consent forms delivered to the participants can be found in the Appendix A.2. The institutional ethics review board has approved the data collection process.

This work contributes to the understanding of LLMs. We do not foresee risk beyond the minimal risk posed by LLM evaluation work. We acknowledge that using LLMs in real-world scenarios could significantly impact human behaviors, raising the need for model transparency, safety, security, and interpretability. We will open-source the synthetic logical reasoning dataset upon publication.

References

01.AI. 2024. Yi: Open Foundation Models by 01.AI. 587

AI@Meta. 2024. The Llama 3 Herd of Models.
Roberta Ballarin. 2023. Modern origins of modal logic. In Edward N. Zalta and Uri Nodelman, editors, <i>The</i> <i>Stanford Encyclopedia of Philosophy</i> , Fall 2023 edi- tion. Metaphysics Research Lab, Stanford University.
Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting Linear Mixed-Effects Models Using Ime4 . <i>Journal of Statistical Software</i> , 67(1).
Catarina G. Belem, Markelle Kelly, Mark Steyvers, Sameer Singh, and Padhraic Smyth. 2024. Percep- tions of Linguistic Uncertainty by Language Models and Humans.
Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2021. Transformers as soft reasoners over language. In <i>IJCAI</i> .
Vittoria Dentella, Fritz Günther, and Evelina Leivada. 2023. Systematic testing of three language mod- els reveals low language accuracy, absence of re- sponse stability, and a yes-response bias. <i>Pro-</i> <i>ceedings of the National Academy of Sciences</i> , 120(51):e2309583120.
Tiwalayo Eisape, Michael Tessler, Ishita Dasgupta, Fei Sha, Sjoerd Steenkiste, and Tal Linzen. 2024. A systematic comparison of syllogistic reasoning in humans and language models. In <i>NAACL</i> .
Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. 2024. Open Ilm leaderboard v2. https://huggingface. co/spaces/open-llm-leaderboard/open_llm_ leaderboard.
Hila Gonen, Srini Iyer, Terra Blevins, Noah Smith, and Luke Zettlemoyer. 2023. Demystifying prompts in language models via perplexity estimation. In <i>Find-</i> <i>ings of ACL: EMNLP</i> .
Christopher Hahn, Frederik Schmitt, Jens U Kreber, Markus Norman Rabe, and Bernd Finkbeiner. 2021. Teaching temporal logics to neural networks. In <i>ICLR</i> .
Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun, Eka- terina Zubova, Yujie Qiao, Matthew Burtell, David Peng, Jonathan Fan, Yixin Liu, Brian Wong, Mal- colm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Shafiq Joty, Alexander R. Fab- bri, Wojciech Kryscinski, Xi Victoria Lin, Caiming Xiong, and Dragomir Radev. 2022. FOLIO: Natural Language Reasoning with First-Order Logic.

- Wesley H. Holliday and Matthew Mandelkern. 2024. Conditional and Modal Reasoning in Large Language Models. ArXiv:2401.17169.
- Jennifer Hu and Roger Levy. 2023. Prompting is not a substitute for probability measurements in large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language

598

599

600

601

602

603

604

605 606

607

608

609 610

611

612

613

614 615

616 617

618 619

620

621

622

623

624

625

626

627

628

629

630 631

632

633

634

635

636

637

638

639

640

641

588

589

- *Processing*, pages 5040–5060, Singapore. Associa-tion for Computational Linguistics.
 - Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *ICML*, pages 9118–9147. PMLR.
 - Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B.

651

653

658

668

670

671

673

674

675

678

679

684

686

692

- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of Experts.
- Philip Nicholas Johnson-Laird. 1983. *Mental models: Towards a cognitive science of language, inference, and consciousness.* 6. Harvard University Press.
- Henry A Kautz, Bart Selman, et al. 1992. Planning as satisfiability. In *ECAI*, volume 92, pages 359–363. Citeseer.
- Saul A. Kripke. 1959. A Completeness Theorem in Modal Logic. *The Journal of Symbolic Logic*, 24(1):1–14.
- Saul A. Kripke. 1963. Semantical Analysis of Modal Logic I Normal Modal Propositional Calculi. *Mathematical Logic Quarterly*, 9(5-6):67–96.
- Andrew K Lampinen, Ishita Dasgupta, Stephanie C Y Chan, Hannah R Sheahan, Antonia Creswell, Dharshan Kumaran, James L McClelland, and Felix Hill. 2024. Language models, like humans, show content effects on reasoning tasks. *PNAS Nexus*, 3(7):pgae233.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. LogiQA: A Challenge Dataset for Machine Reading Comprehension with Logical Reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, pages 3622–3628, Yokohama, Japan. International Joint Conferences on Artificial Intelligence Organization.

- R. Thomas McCoy, Shunyu Yao, Dan Friedman, Mathew D. Hardy, and Thomas L. Griffiths. 2024. Embers of autoregression show how large language models are shaped by the problem they are trained to solve. *Proceedings of the National Academy of Sciences*, 121(41):e2322420121.
- Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2023. Locally typical sampling. *Transactions of the Association for Computational Linguistics*, 11:102–121.
- Microsoft. 2023. Phi-2: The surprising power of small language models. *Microsoft Research Blog*.
- Microsoft. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone.
- Kanishka Misra and Najoung Kim. 2024. Generating novel experimental hypotheses from language models: A case study on cross-dative generalization. *arXiv preprint arXiv:2408.05086*.
- Santiago Ontanon, Joshua Ainslie, Vaclav Cvicek, and Zachary Fisher. 2022. LogicInference: A new Datasaet for Teaching Logical Inference to seq2seq Models. In *ICLR2022 Workshop on the Elements of Reasoning: Objects, Structure and Causality.*
- OpenAI. 2024. Learning to reason with llms. https://openai.com/index/ learning-to-reason-with-llms/.
- 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, Nisarg Patel, Neeraj Varshney, Mutsumi Nakamura, Man Luo, Santosh Mashetty, Arindam Mitra, and Chitta Baral. 2024. LogicBench: Towards systematic evaluation of logical reasoning ability of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers).*
- Marco Ragni, Hannah Dames, Daniel Brand, and Nicolas Riesterer. 2019. When Does a Reasoner Respond: Nothing Follows?: 41st Annual Meeting of the Cognitive Science Society. *Proceedings of the 41st Annual Conference of the Cognitive Science Society*, pages 2640–2645.
- Stephen W Raudenbush. 2002. Hierarchical linear models: Applications and data analysis methods. Advanced Quantitative Techniques in the Social Sciences Series/SAGE.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez,

752

- 765 766 767 768
- 76 76 77
- 771 772 773

774

775

- 776 777 778 778
- 779 780
- 781
- 7 7
- 785 786

78

788 789

790

- 791 792
- 793 794

795

796 797 798

7 8

800 801 802 Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code Llama: Open Foundation Models for Code.

- Abulhair Saparov and He He. 2022. Language Models Are Greedy Reasoners: A Systematic Formal Analysis of Chain-of-Thought. In *The Eleventh International Conference on Learning Representations*.
- Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi, Najoung Kim, and He He. 2023. Testing the General Deductive Reasoning Capacity of Large Language Models Using OOD Examples. *Advances in Neural Information Processing Systems*, 36:3083–3105.
- Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I Wang. 2022. Natural language to code translation with execution. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*.
- Stuart M. Shieber. 1993. The problem of logical form equivalence. *Computational Linguistics*, 19(1):179– 190.
- Damien Sileo and Antoine Lernould. 2023. MindGames: Targeting theory of mind in large language models with dynamic epistemic modal logic. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4570–4577.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2022. Entailer: Answering questions with faithful and truthful chains of reasoning. In *EMNLP*.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Hugo Touvron, Louis Martin, and Kevin Stone. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models.
- Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-tse Huang, Pinjia He, Wenxiang Jiao, and Michael R. Lyu. 2024. A & B == B & A: Triggering Logical Reasoning Failures in Large Language Models.
- Yimei Xiang. 2019. Two types of higher-order readings of wh-questions. *Proceedings of the 22nd Amsterdam Colloquium.*
- Shi Zong and Jimmy Lin. 2024. Categorical syllogisms revisited: A review of the logical reasoning abilities of llms for analyzing categorical syllogism. *arXiv preprint arXiv*:2406.18762.

A Additional Experiment Details

A.1 LLM Experiment Details

All LLMs used are obtained from Hugging Face checkpoints. Time and compute power requirements vary, the largest llama-3-70b model takes around 2 hours on NVIDIA A6000 GPU to obtain all results in §4.

803

804

805

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

A.2 Human Experiment Details

Participant consent. We use the following language to obtain consent from participants, where our institution name is replaced with *the Anonymous Institution* to protect the anonymity of submission.

This study is part of a scientific research project at the Anonymous Institution. Your decision to complete this study is voluntary. There is no way for us to identify you. The only information we will have, in addition to your responses, is the demographic information you provided to Prolific and the time at which you completed the survey. The results of the research may be presented at scientific meetings or published in scientific journals. Clicking on the button below indicates that you are at least 18 years of age and agree to complete this study voluntarily. Press the button below to start the experiment.

Participant instructions. We use keys F and J, which are roughly symmetric on a standard English keyboard, to collect participant responses. Half of the participants see the following instruction:

In this study, you will be presented with two statements followed by a question. Your task is to answer either Yes or No to the question, based on the information provided in the statements. Please respond quickly and accurately by pressing "F" for Yes, and "J" for No.

To mitigate the possible bias introduced by the dominant hand, we have the other half of the participants see instruction with reversed keys: *In this study, you will be presented with two statements followed by a question. Your task is to answer either Yes or No to the question, based on the information provided in the statements. Please respond quickly and accurately by pressing "F" for No, and "J" for Yes.*

Participant wage. We offer participants an hourly wage of 1.5 times Prolific's minimum wage. The duration is determined by the median completion time among all participants.

852

871

874

876

877

885

B Extra Details of the Dataset

B.1 Data Synthesis Pipeline

To better explain the data synthesis process, we provide a detailed visualization of our pipeline in Figure A1.

B.2 Considerations in Translating Logical Form to Natural Language

During the interpretation process, another key point is to assign independent interpretations to variables. Deciding the dependency also involves common sense knowledge. For example, consider the premises $\neg p \rightarrow q$ and q. If we interpret p :="Jane is inside the house" and q := "Jane is out" to proposition variables p and q, the two variables are possibly not independent. According to common sense, "Jane is not inside the house" ($\neg p$) correlates with or is even equivalent to "Jane is out" (q). Logically, $\{\neg p \rightarrow q, q\} \nvDash \neg p$; however, with the extra premise $\neg p \leftrightarrow q$ given by common sense, people may conclude that $\neg p$.¹⁰

Besides, natural language is ambiguous—one sentence in natural language can come from multiple logical forms under the same interpretation. We use present tense and progressive aspect to encourage a reading of imaginary ongoing events, corresponding to the alethic modality. Such events are less likely to induce LLM's or human's individual bias, as they are unrelated to factual knowledge or moral judgements. Also, we always use two full verb phrases, ruling out sentences like "*Jane is eating apples or oranges*," so the two events are less likely to be mutually exclusive. In this way, we can reduce the ambiguity of the questions in our dataset.

B.3 Data Samples

All logic forms and corresponding natural language sentences can be found in Table A1.

The exact prompt format is as follows:

Consider the following statements:\n Jane is watching a show or John is reading a book.\n Jane isn't watching a show.\n Question: Based on these statements, can we infer that John is reading a book?\n Answer:<eof>

C Additional Experiments

C.1 Extra Experiment: Introduction Rule of Modality

We report the results on the necessitation rule and 892 its variants here, as these rules are obscure and 893 verbose to be articulated in natural language: 894

$$\{\varphi\} \vdash \Box \varphi,$$
 (necessitation rule) 895

$$\{\varphi\} \vdash \Diamond \varphi,$$
 896

890

891

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

$$\{\varphi\} \vdash \varphi.$$
 89

Its natural language form is as follows:

Jane is watching a show.

- (□) Can we infer that it's certain that Jane is watching a show?
- (\diamond) Can we infer that it's possible that Jane is watching a show?
- (\emptyset) Can we infer that Jane is watching a show?

All three variants are paired with 1000 logic interpretations. As they are all rules of inference, the ground truth answer is always Yes. Overall accuracy is shown in Table A2, where across all LLMs, the necessitation rule has the lowest accuracy. This echoes the necessity modality's tendency to be rejected discussed in §4.2.3.

We further fit a linear mixed-effects model similar to Eq. (4), except that the argument form effect is now constant across all data points. The mixedeffects model yields a marginal R^2 of 0.391 and a conditional R^2 of 0.745. Estimated marginal means shows that the accuracy on \emptyset is 0.171 less than \diamond , but 0.371 higher than \Box , with both differences significant at p < 0.0001. This further suggests that modality serves as an important factor on logic reasoning performance.

C.2 Extra Experiment: Distribution of Modalities

Besides the necessitation rule, *distribution axiom* is the other fundamental axiom in normal modal logic. It can be transformed into the rule shown in Eq. (A1), and plugging in the definition of \lor in Eq. (2) gives the rule shown in Eq. (A2). Notice that Eq. (A2) closely resembles rule \lor^{L} 's variant with necessity, as shown in Eq. (A3), except the different scope of the necessity operator and the position of the negation operator. Moving the negation operator out of the necessity operator will result in

¹⁰This confounding factor affects the examples in Figure 10 of Holliday and Mandelkern (2024).

Validity	Modality	Argument Form	Logical Form	Natural Language
F	Ø	\vee^{L}	$\{p \lor q, \neg p\} \vdash q$	Jane is watching a show or John is reading a book. Jane isn't watching a show.
	Ø	\vee^{R}	$\{p \lor q, \neg q\} \vdash p$	Can we infer that John is <u>reading a book</u> ? Jane is <u>watching a show or John is reading a book</u> . John isn't <u>reading a book</u> . Can we infer that Jane is watching a show?
	Ø	\rightarrow^{L}	$\{\neg p \to q, \neg p\} \vdash q$	If Jane isn't watching a show, then John is reading a book. Jane isn't watching a show. Can we infer that John is reading a book?
	Ø	\rightarrow^{R}	$\{\neg p \to q, \neg q\} \vdash p$	If Jane isn't watching a show, then John is reading a book. John isn't reading a book. Can we infer that Jane is watching a show?
		\vee^{L}	$\{\Box p \lor \Box q, \neg \Box p\} \vdash \Box q$	It's certain that Jane is watching a show or it's certain that John is reading a book. It's uncertain whether Jane is watching a show. Can we infer that it's certain that John is reading a book?
		\vee^{R}	$\{\Box p \lor \Box q, \neg \Box q\} \vdash \Box p$	It's certain that Jane is watching a show or it's certain that John is reading a book. It's uncertain whether John is reading a book.
		\rightarrow^{L}	$\{\neg \Box p \to \Box q, \neg \Box p\} \vdash \Box q$	a book. It's uncertain whether Jane is watching a show, then it's certain that John is reading a book. It's uncertain whether Jane is watching a show.
		$ ightarrow^{ m R}$	$\{\neg \Box p \to \Box q, \neg \Box q\} \vdash \Box p$	If it's uncertain whether Jane is watching a show, then it's certain that John is reading a book. It's uncertain whether Jane is watching a show, then it's certain that John is reading a book.
	\diamond	\vee^{L}	$\{\Diamond p \lor \Diamond q, \neg \Diamond p\} \vdash \Diamond q$	Can we infer that it's certain that <u>Jane</u> is <u>watching a show</u> ? It's possible that <u>Jane</u> is <u>watching a show</u> or it's possible that <u>John</u> is <u>reading a book</u> . It's impossible that <u>Jane</u> is <u>watching a show</u> .
	\diamond	\vee^{R}	$\{\Diamond p \lor \Diamond q, \neg \Diamond q\} \vdash \Diamond p$	Can we infer that it's possible that John is reading a book? It's possible that John is reading a show or it's possible that John is reading a book. It's impossible that John is reading a book.
	\diamond	\rightarrow^{L}	$\{\neg \Diamond p \to \Diamond q, \neg \Diamond p\} \vdash \Diamond q$	If it's impossible that Jane is watching a show, then it's possible that John is reading a book. It's impossible that Jane is watching a show.
	\$	\rightarrow^{R}	$\{\neg \Diamond p \to \Diamond q, \neg \Diamond q\} \vdash \Diamond p$	If it's impossible that John is reading a book. It's impossible that John is reading a book.
¥	Ø	\vee^{L}	$\{p \lor q,q\} \nvdash \neg p$	Jane is watching a show or John is reading a book. John is reading a book.
	Ø	$\vee^{\mathbf{R}}$	$\{p \lor q, p\} \nvDash \neg q$	Jane is watching a show or John is reading a book.
	Ø	\rightarrow^{L}	$\{\neg p \to q,q\} \nvDash \neg p$	If Jane isn't watching a show, then John is reading a book. John is reading a book. Can we infer that Jane isn't watching a show?
	ø	$ ightarrow^{ m R}$	$\{\neg p \to q, p\} \nvDash \neg q$	If Jane isn't watching a show, then John is reading a book. Jane is watching a show.
		\vee^{L}	$\{\Box p \lor \Box q, \Box q\} \nvdash \neg \Box p$	It's certain that Jane is watching a show or it's certain that John is reading a book. It's certain that John is reading a book. Can we infer that it's uncertain whether Jane is watching a show?
		\vee^{R}	$\{\Box p \lor \Box q, \Box p\} \nvdash \neg \Box q$	It's certain that Jane is watching a show or it's certain that Jane is watching a show. It's certain that Jane is watching a show. Can we infer that it's uncertain whether John is reading a book?
		\rightarrow^{L}	$\{\neg\Box p\to\Box q,\Box q\}\nvDash\neg\Box p$	If it's uncertain whether Jane is watching a show, then it's certain that John is reading a book. It's certain that John is reading a book. Con we infer that if's uncertain whether Jane is watching a show?
		\rightarrow^{R}	$\{\neg \Box p \to \Box q, \Box p\} \nvdash \neg \Box q$	If it's uncertain whether Jane is watching a show, then it's certain that John is reading a book. It's certain that Jane is watching a show.
	\diamond	\vee^{L}	$\{\Diamond p \lor \Diamond q, \Diamond q\} \nvdash \neg \Diamond p$	Can we infer that it's uncertain whether John is reading a book? It's possible that Jane is watching a show or it's possible that John is reading a book. It's possible that John is reading a book.
	\diamond	\vee^{R}	$\{\Diamond p \lor \Diamond q, \Diamond p\} \nvdash \neg \Diamond q$	Can we infer that it's impossible that Jane is watching a show? It's possible that Jane is watching a show or it's possible that John is reading a book. It's possible that Jane is watching a show.
	\$	$ ightarrow^{ m L}$	$\{\neg \Diamond p \to \Diamond q, \Diamond q\} \nvDash \neg \Diamond p$	If it's impossible that John is reading a book. It's possible that John is reading a book. It's possible that John is reading a book.
	\diamond	$ ightarrow^{ m R}$	$\{\neg \Diamond p \to \Diamond q, \Diamond p\} \nvdash \neg \Diamond q$	Gan we infer that it's impossible that <u>Jane</u> is <u>watching a show</u> ? If it's impossible that <u>Jane</u> is <u>watching a show</u> , then it's possible that <u>John</u> is <u>reading a book</u> . It's possible that <u>Jane is watching a show</u> . Can we infer that it's impossible that <u>John</u> is <u>reading a book</u> ?

Table A1: Samples of all logical forms and corresponding natural language sentences.



Figure A1: The data synthesis pipeline: for each variable in logic forms (§3.1) we assign meanings to them to obtain the natural language question-answering pairs (§3.2).

	Ø		\diamond
mistral-7b	0.998	0.885	0.999
mistral-8x7b	0.957	0.540	0.987
llama-2-7b	0.768	0.013	0.920
llama-2-13b	0.368	0.004	0.829
llama-2-70b	0.511	0.051	0.834
llama-3-8b	0.398	0.225	0.783
llama-3-70b	0.674	0.384	0.794
yi-34b	0.960	0.382	0.999
phi-2	0.814	0.226	0.892
phi-3-mini	0.992	0.925	0.994

Table A2: Overall accuracy of the necessitation rule and its modality variants on each model.

a fallacy (Eq. A4).

928

929

930

931

932

933

936

937

$$\{\Box(\varphi \to \psi), \Box\varphi\} \vdash \Box\psi, \tag{A1}$$

$$\{\Box(\varphi \lor \psi), \Box \neg \varphi\} \vdash \Box \psi, \tag{A2}$$

$$\{\Box\varphi \lor \Box\psi, \neg \Box\varphi\} \vdash \Box\psi, \tag{A3}$$

$$\{\Box(\varphi \lor \psi), \neg \Box \varphi\} \nvDash \Box \psi. \tag{A4}$$

We say (A2) to (A4) are of argument form theorem, base and spurious, respectively. See Table A3 for the logical forms and their ground truth we used to study the distribution of modalities. The natural language form is as follows:

	It's certain that if Freddy is not going
	shopping, then Coy is making dinner.
(theorem)	It's certain that Freddy is not going shop-
	ping.
(spurious)	It's uncertain whether Freddy is going
	shopping.
	Can we infer that it's certain that Coy is
	making dinner?

938This group of rules and fallacies comes from939the fact that the necessity modality \Box is not dis-940tributive to disjunction, i.e. $\Box(\varphi \lor \psi) \nvDash \Box \varphi \lor \Box \psi$ 941(Xiang, 2019, Ex. 5). In contrast, the possibility

Modality	Argument Form	Logical Form
Ø	base	$\varphi \lor \psi, \neg \varphi \vdash \psi$
	base	$\Box \varphi \lor \Box \psi, \neg \Box \varphi \vdash \Box \psi$
	theorem	$\Box(\varphi \lor \psi), \Box \neg \varphi \vdash \Box \psi$
	spurious	$\Box(\varphi \lor \psi), \neg \Box \varphi \nvDash \Box \psi$
\diamond	base	$\Diamond \varphi \lor \Diamond \psi, \neg \Diamond \varphi \vdash \Diamond \psi$
\diamond	theorem	$\Diamond(\varphi \lor \psi), \Diamond \neg \varphi \vdash \Diamond \psi$
\diamond	spurious	$\Diamond(\varphi \lor \psi), \neg \Diamond \varphi \vdash \Diamond \psi$

Table A3: Logical forms and their ground truth to study the distribution of modalities. Only the spurious form of the necessity modality (marked by <u>underline</u>) has a ground truth of false.

modality \diamondsuit is distributive to disjunction. This particular case could have served as a material to test the LLM's knowledge of the asymmetry between the two modalities, yet in §4.2.3 we showed that there is a bias towards rejection on the necessity modality. As the false case of the disjunction is on the necessity modality, this bias confounds the experiment. 942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

We fit a linear mixed-effects model similar to Eq. (4) to the data,

$$Acc_{soft} \sim Modality \times ArgForm + Perplexity + (1 + Perplexity | LLM),$$

with an interaction term between the modality and argument form. On the theorem form compared to the base form, the necessity modality \Box has a 0.173 higher estimated marginal means with p <0.0001 significance, yet the possibility modality \diamond has a 0.071 lower estimated marginal means. On the spurious form compared to the base form, the \Box has a 0.312 higher means, and the \diamond has no significant difference. On both forms, $\diamond \succ \Box$ in terms of accuracy still holds at a slight margin of 0.110 and 0.047 respectively.

965	To verify whether on \Box the performance increase
966	on spurious form is due to the rejection bias, we
967	fit a linear mixed-effects model with the relative
968	probability of answering Yes as dependent variable.
969	Results show that on spurious form compared to
970	the base form, the effect of \Box 's tendency to answer
971	Yes is only 0.060 lower, indicating the rejection
972	bias of the base form is still present. Therefore,
973	we hypothesize that the LLM's performance on
974	recognizing the fallacy of necessity distribution
975	over disjunction is hindered by the rejection bias
976	on the necessity modality.