

# BENCHMARKING LOGICAL REASONING INCONSISTENCIES IN LOCAL LARGE LANGUAGE MODELS: EVIDENCE FROM MULTI-DOMAIN EVALUATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We present systematic evidence of logical reasoning limitations in local large language models through MREB (Multimodal Reasoning and Ethics Benchmark), focusing on deduction, induction, and consistency across related questions. Our evaluation of four prominent local models reveals significant logical reasoning deficits, with performance ranging from 48-60% on logical tasks while achieving 84-92% on ethics questions that require similar reasoning patterns. We identify three critical failure modes: (1) inconsistent logical deduction across semantically equivalent problems, (2) failure to maintain logical consistency when reasoning about related scenarios, and (3) systematic bias toward pattern matching over genuine logical inference. Our findings demonstrate that current local LLMs exhibit fundamental logical reasoning limitations that are masked by strong performance in other cognitive domains, highlighting the need for targeted logical reasoning improvements and more rigorous consistency evaluation frameworks. paragraph.

## 1 INTRODUCTION

Large language models have demonstrated remarkable capabilities across diverse tasks, yet their logical reasoning abilities remain inconsistent and unreliable. While models can achieve high performance on certain reasoning benchmarks, they often fail when logical problems require systematic deduction, consistent application of logical rules, or maintaining coherence across related questions. The deployment of LLMs in local environments introduces additional constraints that may exacerbate these limitations. Resource constraints, quantization effects, and deployment optimizations can impact the precise reasoning capabilities that logical tasks demand. Understanding these limitations is crucial for developing more reliable reasoning systems and establishing appropriate boundaries for LLM deployment in logic-critical applications.

### 1.1 CONTRIBUTIONS

This work provides empirical evidence of logical reasoning limitations in local LLMs through:

1. **Systematic Evaluation:** Comprehensive assessment of logical reasoning capabilities across four local models using 225 carefully designed logical reasoning tasks
2. **Consistency Analysis:** Investigation of logical consistency across related questions and scenarios
3. **Failure Mode Identification:** Detailed analysis of specific logical reasoning failure patterns
4. **Cross-Domain Comparison:** Comparison of logical reasoning performance with other cognitive domains to identify systematic biases

## 2 METHODOLOGY

Our evaluation is designed around two central questions: whether local LLMs can perform structured logical inference at a level suitable for real-world deployment, and whether their reasoning

Table 1: Performance comparison across cognitive domains

Model	Logical Reasoning	Ethics	Coding	Multimodal
llama3.2-vision	60.0%	88.0%	60.0%	84.0%
gemma3	60.0%	84.0%	64.0%	80.0%
llava-llama3	48.0%	84.0%	40.0%	84.0%
llava	48.0%	88.0%	20.0%	92.0%

remains consistent when the same problem is posed in different forms. To answer these, we constructed a benchmark of 225 logical reasoning tasks, developed a consistency evaluation framework, and selected four representative models from the local deployment ecosystem.

**Logical Reasoning Task Design:** Our logical reasoning evaluation consists of 225 tasks spanning three core logical reasoning types. **Deductive Reasoning** (70 tasks) covers classical syllogistic reasoning, modus ponens/tollens, and categorical logic problems requiring precise logical deduction from given premises. **Inductive Reasoning** (75 tasks) encompasses pattern recognition, generalization from examples, and probabilistic reasoning tasks that require drawing general conclusions from specific instances. **Abductive Reasoning** (80 tasks) addresses best explanation selection, causal inference, and hypothesis generation tasks requiring reasoning from observations to likely explanations.

**Consistency Evaluation Framework:** To assess logical consistency, we designed task pairs that test the same logical principles in different contexts. **Semantic Equivalence** tasks present logically identical problems with different surface features. **Contrapositive Consistency** tests whether models maintain logical equivalence between statements and their contrapositives. **Transitivity Preservation** evaluates consistency in transitive reasoning across related question sequences.

**Model Evaluation:** We evaluated four prominent local models from the Ollama ecosystem: **llama3.2-vision:latest** (Meta’s multimodal model), **gemma3:latest** (Google’s Gemma architecture), **llava-llama3:latest** (LLaVA with Llama3 base), and **llava:latest** (the original LLaVA model). All evaluations were conducted under identical local deployment conditions using structured JSON output formatting to ensure consistent response parsing.

### 3 RESULTS

Across all four models and all cognitive domains tested, logical reasoning consistently emerged as the weakest area. The results reveal not only low absolute accuracy but also systematic inconsistency in how models handle logically equivalent problems, suggesting the limitations are structural rather than incidental.

#### 3.1 OVERALL LOGICAL REASONING PERFORMANCE

Table 1 compares model performance across cognitive domains. Logical reasoning scores lagged considerably behind ethics and multimodal tasks despite ethics questions relying on similar underlying reasoning patterns such as conditional logic and consequence evaluation. Models consistently performed 24-40 percentage points better on ethics tasks than on logical reasoning tasks, pointing to a reliance on surface-level cues rather than genuine deductive capability.

#### 3.2 REASONING TYPE ANALYSIS

Performance varied significantly across the three logical reasoning types. **Deductive Reasoning** achieved 45-65% accuracy across models, with best performance on simple syllogisms but significant degradation as premise count increased and frequent errors in contrapositive reasoning. **Inductive Reasoning** achieved 40-55% accuracy, with inconsistent pattern recognition, a tendency to overgeneralize from limited examples, and failure to consider alternative explanations. **Abductive Reasoning** achieved 50-60% accuracy, with models showing a preference for obvious explanations over optimal ones and inconsistent causal reasoning when competing hypotheses were present.

### 3.3 LOGICAL CONSISTENCY FAILURES

Beyond raw accuracy, a more revealing picture emerges from consistency analysis. Models frequently gave contradictory answers to problems that are logically identical, indicating that their responses are shaped more by surface presentation than by underlying logical structure.

#### 3.3.1 SEMANTIC EQUIVALENCE INCONSISTENCY

Models showed a 34% inconsistency rate when presented with logically equivalent problems in different surface forms. For example, Task A (“All cats are mammals. Fluffy is a cat. Is Fluffy a mammal?”) and Task B (“Every feline is a mammal. Whiskers is a feline. Is Whiskers a mammal?”) are logically identical, yet 3 of 4 models gave inconsistent answers across the pair.

#### 3.3.2 CONTRAPOSITIVE REASONING FAILURES

Models demonstrated systematic failure in contrapositive reasoning with a 67% error rate. Given the statement “If it rains, then the ground gets wet,” the correct contrapositive is “If the ground is not wet, then it did not rain.” Models commonly confused this with the converse or inverse statements instead.

#### 3.3.3 TRANSITIVITY VIOLATIONS

Models violated logical transitivity in 45% of relevant task sequences. A representative failure involves a three-step chain where models correctly answer “ $A > B$ ” and “ $B > C$ ” independently, but then incorrectly conclude that “ $A > C$ ” is false, which is a direct contradiction of their prior responses.

## 4 ANALYSIS OF FAILURE MODES

The patterns observed in Section 3 point toward a common underlying cause: models appear to recognize the shape of logical problems rather than reason through them. Three specific mechanisms help explain this behavior.

**Pattern Matching vs. Logical Reasoning:** High performance on ethics questions, which often have obvious socially expected answers, versus poor performance on logical reasoning questions that require systematic deduction suggests models rely on pattern matching rather than genuine logical inference. When logical reasoning questions were reformulated with obvious answer patterns, performance improved by 23% on average, confirming that surface-level recognition rather than deep logical understanding is driving responses.

**Context Length and Complexity Effects:** Performance degraded linearly with the number of premises: 65% accuracy for 2-premise problems, 52% for 3-premise problems, and 38% for problems with 4 or more premises. A similar trend holds for logical depth: single-step inference achieved 68% accuracy, two-step inference 45%, and three-step inference only 31%. Together, these results suggest that models struggle to maintain and integrate logical state as problem complexity grows.

**Local Deployment Impact:** Resource constraints in local deployment likely amplify these weaknesses. Quantization reduces arithmetic precision in ways that may destabilize subtle logical computations. Limited context windows may prevent models from holding all relevant premises in view simultaneously. And the inference speed pressures common in local deployment may further push models toward fast pattern matching rather than deliberate step-by-step deduction.

## 5 IMPLICATIONS FOR LOGICAL REASONING RESEARCH

These findings carry practical consequences for how the community evaluates and develops logical reasoning in language models. The gap between benchmark performance and the consistency failures documented here suggests that existing evaluation practices may be painting an overly optimistic picture.

**Benchmarking Limitations:** Current logical reasoning benchmarks may overestimate model capabilities by using tasks with obvious answer patterns, failing to test consistency across related problems, and not accounting for deployment constraint effects. A model that answers isolated logical questions at 60% accuracy while failing 45% of transitivity chains is considerably less reliable than headline accuracy figures suggest.

**Need for Consistency-Focused Evaluation:** Our results highlight the critical need for evaluation frameworks that explicitly test logical consistency across semantically equivalent problems, maintenance of logical relationships across question sequences, and robustness of logical reasoning under different surface presentations. Accuracy on individual tasks is a necessary but insufficient measure of genuine logical competence.

**Local vs. Cloud Deployment Considerations:** Local deployment constraints may fundamentally limit logical reasoning capabilities in ways that cloud-based evaluations do not capture. This suggests a need for deployment-specific logical reasoning benchmarks, better understanding of how quantization and context window limits affect reasoning chains, and explicit trade-off analyses between deployment efficiency and logical reasoning accuracy.

## 6 FUTURE WORK

The failure modes identified in this study point toward several concrete directions for improving both evaluation methodology and model capability.

**Enhanced Consistency Evaluation:** Future benchmarks should test logical reasoning robustness across multiple surface representations of identical problems, extended reasoning chains that require maintaining logical state over many steps, and cross-domain transfer where the same logical structure appears in different subject areas. Without such tests, evaluation results will continue to overstate true reasoning reliability.

**Targeted Logical Reasoning Improvements:** The specific failure patterns observed like degradation with premise count, contrapositive confusion, and transitivity violations, suggest targeted training interventions. Improving multi-premise reasoning, strengthening understanding of logical equivalence, and reducing reliance on pattern matching in favor of systematic deduction are the most pressing areas for development.

**Deployment-Aware Logical Reasoning:** Future work should investigate how quantization, context limits, and inference constraints interact with logical reasoning performance. This includes developing resource-efficient reasoning methods that do not sacrifice logical consistency, building deployment-specific evaluation suites, and producing trade-off analyses that help practitioners make informed decisions about model selection for reasoning-critical applications.

## 7 CONCLUSION

Our systematic evaluation reveals fundamental limitations in local LLM logical reasoning capabilities, with models achieving only 48-60% accuracy on logical reasoning tasks while performing significantly better on superficially similar ethics questions. More troubling than the accuracy gap is the consistency gap: models routinely contradict themselves across semantically equivalent problems, violate transitivity, and fail at contrapositive reasoning at rates that would render them unreliable in any logic-critical application. The identification of these failure modes and their connection to pattern matching over genuine inference provides concrete targets for improvement. For the logical reasoning research community, these results make a strong case for moving beyond isolated task accuracy toward evaluation frameworks that measure consistency, robustness, and deployment-aware performance. Only through such rigorous evaluation can we develop language models with the systematic logical reasoning capabilities that real-world deployment demands.

## REFERENCES

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. From ‘f’ to ‘a’ on the ny regents science exams: An overview of the aristo project. *AI Magazine*, 41(4):39–53, 2020.

216 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Peter  
217 West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, et al. Faith and fate: Limits of  
218 transformers on compositionality. In *Advances in Neural Information Processing Systems*, vol-  
219 ume 36, 2023.

220 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
221 Steinhardt. Measuring massive multitask language understanding. In *International Conference  
222 on Learning Representations*, 2020.

223 Jiayu Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Evaluating the logical reasoning  
224 ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*, 2023.

225 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John  
226 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models:  
227 Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.

228 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
229 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adria Garriga-Alonso, et al. Beyond the  
230 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint  
231 arXiv:2206.04615*, 2022.

232 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, Denny  
233 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances  
234 in Neural Information Processing Systems*, volume 35, pp. 24824–24837, 2022.

## 235 A RELATED WORK CONTEXT AND CITATION RATIONALE

236 This appendix provides additional context for our key citations and their relevance to logical reason-  
237 ing evaluation in large language models.

### 238 A.1 CAPABILITY SCALING AND EMERGENCE

239 **Srivastava et al. (2022) - BIG-bench:** The Beyond the Imitation Game benchmark (Srivastava  
240 et al., 2022) established that reasoning abilities emerge as emergent properties of scale in large  
241 language models. Their work demonstrated that logical reasoning performance follows predictable  
242 scaling curves, with qualitative improvements appearing at specific model sizes. This foundational  
243 work supports our observation that current local models may exhibit reasoning limitations due to  
244 scale constraints inherent in local deployment scenarios.

245 **Rae et al. (2021) - Gopher:** The Gopher scaling study (Rae et al., 2021) provided crucial insights  
246 into when reasoning capabilities begin to manifest during model scaling. Their analysis of reason-  
247 ing emergence across different model sizes contextualizes our findings about local model limitations,  
248 suggesting that resource-constrained deployment may fundamentally limit access to reasoning ca-  
249 pabilities that only emerge at larger scales.

250 **Relevance to Our Work:** These scaling studies establish the theoretical foundation for understand-  
251 ing why local models (which are often smaller or quantized versions) may exhibit the logical rea-  
252 soning deficits we observe. The 48-60% performance range in our evaluation may reflect models  
253 operating below the scale threshold where robust logical reasoning emerges.

### 254 A.2 PROMPT-BASED REASONING ENHANCEMENT

255 **Wei et al. (2022) - Chain-of-Thought:** The seminal chain-of-thought prompting work (Wei et al.,  
256 2022) demonstrated that explicit reasoning traces can significantly improve logical reasoning perfor-  
257 mance. This work is directly relevant to our consistency evaluation framework—the failure modes  
258 we identify (semantic equivalence inconsistency, contrapositive reasoning failures) suggest that cur-  
259 rent models lack the systematic reasoning processes that chain-of-thought prompting attempts to  
260 induce.

261 **Implications for Our Findings:** The success of chain-of-thought prompting in improving reason-  
262 ing performance supports our hypothesis that current models rely on pattern matching rather than  
263  
264  
265  
266  
267  
268  
269

genuine logical inference. Our diagnostic finding that performance improved 23% when logical reasoning questions were reformulated with obvious answer patterns aligns with Wei et al.’s observation that explicit reasoning structure enhances performance.

### A.3 LOGICAL AND STRUCTURED REASONING LIMITATIONS

**Liu et al. (2023) - ChatGPT & GPT-4 Reasoning Evaluation:** Liu et al.’s systematic evaluation of logical reasoning in state-of-the-art models (Liu et al., 2023) provides crucial context for our local model findings. Their identification of logical reasoning failures in even the most capable cloud-based models suggests that the limitations we observe in local models represent broader challenges in LLM logical reasoning rather than deployment-specific issues alone.

**Dziri et al. (2023) - Compositionality Limits:** The ”Faith and Fate” study (Dziri et al., 2023) established fundamental limits in transformer compositionality, directly supporting our findings about multi-premise reasoning degradation. Their work on compositional generalization failures provides theoretical grounding for our observed performance degradation with increased logical complexity (65%  $\rightarrow$  52%  $\rightarrow$  38% accuracy as premises increase).

**Connection to Our Consistency Analysis:** Both studies support our core finding that logical reasoning pattern matching. The consistency failures we document (34% semantic equivalence inconsistency, 67% contrapositive reasoning errors) align with their broader findings about LLM reasoning limitations.

### A.4 BENCHMARKING AND MULTI-DOMAIN EVALUATION

**Hendrycks et al. (2020) - MMLU:** The Massive Multitask Language Understanding benchmark (Hendrycks et al., 2020) established the framework for broad academic evaluation that inspired our multi-domain approach. However, MMLU’s focus on knowledge recall rather than reasoning consistency highlights the gap our work addresses like the need for evaluation frameworks that test logical coherence across related questions.

**Clark et al. (2020) - Aristo Project:** The Aristo science reasoning work (Clark et al., 2020) demonstrated the challenges of systematic reasoning in academic domains. Their progression from ’F’ to ’A’ performance on science exams provides a model for the type of systematic reasoning improvement needed to address the logical consistency failures we identify.

**Methodological Influence:** Both benchmarking efforts informed our decision to include cross-domain comparison in MREB. The significant performance gap we observe between logical reasoning (48-60%) and ethics tasks (84-92%) builds on their multi-domain evaluation approaches while revealing domain-specific reasoning limitations.

### A.5 SYNTHESIS AND FUTURE DIRECTIONS

The convergence of evidence from scaling studies (Srivastava et al., 2022; Rae et al., 2021), reasoning enhancement techniques (Wei et al., 2022), limitation analyses (Liu et al., 2023; Dziri et al., 2023), and benchmarking frameworks (Hendrycks et al., 2020; Clark et al., 2020) supports our central thesis: current LLMs exhibit fundamental logical reasoning limitations that require targeted evaluation and improvement approaches.

Our contribution extends this body of work by:

1. Providing systematic evidence of reasoning limitations in resource-constrained deployment scenarios
2. Introducing consistency-focused evaluation methodology that reveals hidden reasoning failures
3. Demonstrating the inadequacy of single-task reasoning evaluation for understanding true logical capabilities

Future work should build on these foundations to develop reasoning enhancement techniques specifically designed for local deployment constraints while maintaining the rigorous consistency evaluation standards our work establishes.