

State-of-the-Art Large Language Model Performance on Reasoning Benchmarks

Assignee Research

June 5, 2026

Abstract

This report synthesises findings from 6 peer-reviewed papers addressing the following research question: What are the state-of-the-art large language model results on reasoning benchmarks published recently v6. 10 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Lion: Adversarial Distillation of Proprietary Large Language Models. Research question: What are the state-of-the-art large language model results on reasoning benchmarks published recently v6.

2 Methodology

Systematic literature search across multiple databases yielded 6 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

3 Results

6 papers retrieved. 10 claims extracted; 0 independently verified. Quality review score: 3.7/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce

errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
Adversarial Knowledge Distillation (AKD) methodologies necessitate accessibility to the weights or gradients of the teacher	×	0.11
The Lion framework consists of three stages in an iteration: an imitation stage, a discrimination stage, and a generative stage	×	0.09
In the Lion framework, the proprietary teacher LLM is prompted to serve as a 'referee' to discriminate hard instructions	×	0.13
In the Lion framework, the proprietary teacher LLM is prompted to serve as a 'generator' to produce new instructions that	×	0.11
The Vicuna-Instructions dataset consists of 80 questions spanning 9 distinct task categories.	×	0.03
Setting1 of the evaluation protocol uses GPT-4 to automatically assess response quality on a scale of 1 to 10 between a	×	0.04
Wang et al. (2023) identified a systematic bias in GPT-4 automatic evaluation and proposed Multiple Evidence Calibration	×	0.03
AGIEval is a benchmark designed to evaluate the capability of foundation models in the context of human-centric standard	×	0.04
Previous knowledge distillation works employ a unidirectional approach where the teacher imparts knowledge to the student	×	0.12
The authors applied the AKD framework to transfer knowledge from ChatGPT to the open-source LLaMA model.	×	0.14

References

- <http://arxiv.org/abs/2305.14497v2>
- <http://arxiv.org/abs/2305.12870v2>

- <http://arxiv.org/abs/2303.13375v2>