

# Llama-0.72 Benchmark Performance Across Reasoning Mathematics Coding and Language Tasks

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: What are the benchmark performance scores of Llama-0.72 on reasoning mathematics coding and language understanding tasks. 8 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: MMLU-Pro: A More Robust and Challenging Multi-Task Language Understanding Benchmark. Research question: What are the benchmark performance scores of Llama-0.72 on reasoning mathematics coding and language understanding tasks.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.5/10.

## 3 Results

16 papers retrieved. 8 claims extracted; 0 independently verified. Quality review score: 3.5/10.

## 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
GPT-4o achieves an overall accuracy of 72.6% on the MMLU-Pro benchmark.	×	0.07
Phi-3-medium-4k-instruct (14B parameters) and Phi-3-mini-4k-instruct (3.8B parameters) perform exceptionally well on the	×	0.04
Llama-3-70B-Instruct achieves an accuracy of 56.2% on the MMLU-Pro benchmark.	×	0.07
GPT-4o scores over 70% accuracy in Math and Physics on the MMLU-Pro benchmark.	×	0.08
Mistral-7B-v0.1 scores just over 20% accuracy in Math and Physics on the MMLU-Pro benchmark.	×	0.08
DeepSeek-V2-Chat underperforms relative to its peers in History and Psychology on the MMLU-Pro benchmark.	×	0.05
Engineering and Law consistently scored lower among the 14 subjects evaluated on the MMLU-Pro benchmark.	×	0.05
GPT-4o’s performance was analyzed through a detailed review of 120 randomly selected erroneous predictions.	×	0.03

## References

- <http://arxiv.org/abs/2406.01574v6>
- <http://arxiv.org/abs/2509.25160v1>
- <http://arxiv.org/abs/2410.12381v3>