

# Grok-Imagine Benchmark Performance Across Reasoning Mathematics Coding and Language Tasks

Assignee Research

June 7, 2026

## Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: What are the benchmark performance scores of Grok-Imagine on reasoning mathematics coding and language understanding tasks. 7 claims were extracted from source literature; 7 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Large Language Models: A Survey. Research question: What are the benchmark performance scores of Grok-Imagine on reasoning mathematics coding and language understanding tasks.

## 2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.3/10.

## 3 Results

9 papers retrieved. 7 claims extracted; 7 independently verified. Quality review score: 8.3/10.

## 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
Large Language Models (LLMs) have drawn a lot of attention due to their strong performance on a wide range of natural la	✓	0.42
LLMs' ability of general-purpose language understanding and generation is acquired by training billions of model's param	✓	0.39
The research area of LLMs, while very recent, is evolving rapidly in many different ways.	✓	0.28
In this paper, we review some of the most prominent LLMs, including three popular LLM families (GPT, LLaMA, PaLM), and d	✓	0.40
We also give an overview of techniques developed to build, and augment LLMs.	✓	0.23
We then survey popular datasets prepared for LLM training, fine-tuning, and evaluation, review widely used LLM evaluatio	✓	0.49
Finally, we conclude the paper by discussing open challenges and future research directions.	✓	0.29

## References

- <https://doi.org/10.48550/arxiv.2402.06196>
- <https://doi.org/10.3389/frdem.2024.1385303>
- <https://doi.org/10.3386/w34202>