

# SOVEREIGN: Holistic evaluation LLM benchmark results MMLU HumanEval GSM8K MATH SWE-bench accuracy scores table 2024

SOVEREIGN Research Kernel  
Autonomous draft — Owner review required before publication

May 29, 2026

## Abstract

Abstract The rapid evolution of large language models (LLMs) has driven a transformative shift in artificial intelligence (AI), reshaping both research paradigms and practical applications. Distinguished from their predecessors by unprecedented scale and advanced capabilities, LLMs necessitate new frameworks for understanding their development, behavior, and societal impact. This survey systematically reviews recent advancements in LLM techniques across four key dimensions: (1) pre-training methodologies, which establish core model capabilities through large-scale self-supervised training, arc

## 1 Introduction

Analysis of: A Survey of Large Language Models. Research goal: Holistic evaluation LLM benchmark results MMLU HumanEval GSM8K MATH SWE-bench accuracy scores table 2024.

## 2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

## 3 Results

12 papers retrieved. 9 claims extracted, 9 verified. Tribunal: 9.0/10 \$\rightarrow\$ APPROVE (revision\_round=0). Policy: AUTO\_APPROVE.

## 4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

## 5 Extracted Claims

Claim	Verified	Confidence
The survey reviews LLM advancements across four key dimensions: pre-training methodologies, post-training techniques, ut	✓	0.23
Pre-training methodologies establish core model capabilities through large-scale self-supervised training, architectural	✓	0.36
Post-training techniques include supervised fine-tuning and reinforcement learning.	✓	0.19
Post-training techniques adapt foundational models to downstream tasks and enhance their alignment and safety.	✓	0.27
Utilization strategies include in-context learning, prompt engineering, and agentic reasoning.	✓	0.20
Utilization strategies optimize real-world deployment and enable effective interaction with external environments.	✓	0.24
Evaluation methods encompass benchmarks for core language capabilities, reasoning, and safety.	✓	0.19
The survey identifies critical research issues concerning theoretical foundations, efficient scaling, alignment, and age	✓	0.26
Large language models (LLMs) are distinguished from their predecessors by unprecedented scale and advanced capabilities.	✓	0.26

## References

- <https://doi.org/10.1007/s11704-026-60308-3>
- <https://doi.org/10.48550/arxiv.2505.09388>
- <https://doi.org/10.48550/arxiv.2412.19437>