

To what extent does zero-shot reasoning accuracy degrade in Llama3, Codestral, and Deepseek R1 when forecastin

Assignee Research

May 29, 2026

Abstract

The field of artificial intelligence has undergone a revolution from foundational Transformer architectures to reasoning-capable systems approaching human-level performance. We present LLMOrbit, a comprehensive circular taxonomy navigating the landscape of large language models spanning 2019-2025. This survey examines over 50 models across 15 organizations through eight interconnected orbital dimensions, documenting architectural innovations, training methodologies, and efficiency patterns defining modern LLMs, generative AI, and agentic systems. We identify three critical crises: (1) data sca

1 Introduction

This paper examines: LLMOrbit: A Circular Taxonomy of Large Language Models -From Scaling Walls to Agentic AI Systems. Research question: To what extent does zero-shot reasoning accuracy degrade in Llama3, Codestral, and Deepseek R1 when forecasting time-series anomalies in domain-shifted datasets compared to in-domain benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 1 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.2/10.

3 Results

1 papers retrieved. 14 claims extracted; 5 independently verified. Quality review score: 5.2/10.

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
The LLMOrbit survey examines over 50 models across 15 organizations.	✓	0.18
The survey covers large language models spanning the years 2019 to 2025.	✓	0.16
Data scarcity is projected to result in 9-27T tokens being depleted by 2026-2028.	×	0.14
LLM training costs have grown from \$3M to over \$300M in a 5-year period.	×	0.05
Energy consumption for LLMs has increased by a factor of 22.	×	0.05
Models o1 and DeepSeek-R1 achieve GPT-4 performance using 10x inference compute via test-time compute strategies.	✓	0.22
Quantization techniques provide 4-8x compression for large language models.	×	0.14
Distributed edge computing offers a 10x cost reduction.	×	0.14
The ORPO training method reduces memory usage by 50%.	×	0.09
Phi-4, a 14B parameter model, matches the performance of larger models.	×	0.12
DeepSeek-R1 achieves a score of 79.8% on the MATH benchmark.	×	0.09
Mixture of Experts (MoE) routing provides an 18x efficiency improvement.	×	0.08
Multi-head Latent Attention enables 8x KV cache compression.	✓	0.18
Multi-head Latent Attention enables GPT-4-level performance at a cost of less than \$0.30 per million tokens.	✓	0.19

References

- <https://openalex.org/W7125352730>