

Emergent Reasoning Capabilities in Transformers at Scale: A Multi-Study Synthesis

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What is the relationship between model scale and emergent reasoning capabilities in transformers v10. 8 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.0/10. 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 is the relationship between model scale and emergent reasoning capabilities in transformers v10.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.0/10.

3 Results

13 papers retrieved. 8 claims extracted; 0 independently verified. Quality review score: 4.0/10.

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.	×	0.06
Phi-3-medium-4k-instruct (14B parameters) and Phi-3-mini-4k-instruct (3.8B parameters) perform exceptionally well on the	×	0.03
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.	×	0.07
Mistral-7B-v0.1 scores just over 20% accuracy in Math and Physics on the MMLU-Pro benchmark.	×	0.07
DeepSeek-V2-Chat underperforms relative to its peers in History and Psychology on the MMLU-Pro benchmark.	×	0.04
Engineering and Law consistently scored lower among the 14 subjects evaluated on the MMLU-Pro benchmark.	×	0.04
GPT-4o’s performance was analyzed through a detailed review of 120 randomly selected erroneous predictions on the MMLU-P	×	0.05

References

- <http://arxiv.org/abs/2406.01574v6>
- <http://arxiv.org/abs/2504.14693v2>
- <http://arxiv.org/abs/2503.20786v1>