

# GPT-4o Benchmark Performance Across Reasoning Mathematics Coding and Language Tasks

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: What are the benchmark performance scores of GPT-4o on reasoning mathematics coding and language understanding tasks. 10 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Challenging BIG-Bench Tasks and Whether Chain-of-Thought Can Solve Them. Research question: What are the benchmark performance scores of GPT-4o on reasoning mathematics coding and language understanding tasks.

## 2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.5/10.

## 3 Results

15 papers retrieved. 10 claims extracted; 3 independently verified. Quality review score: 5.5/10.

## 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
PaLM, InstructGPT, and Codex models were evaluated on BBH for answer-only and CoT prompting approaches.	×	0.08
Answer-only prompting typically underestimates language model performance on challenging tasks requiring multiple reason	×	0.08
In the BIG-Bench paper, none of the evaluated models, including PaLM 540B, outperformed human-rater baselines on any of	×	0.15
PaLM 540B with answer-only prompting outperforms the average human-rater on 6 out of 23 BBH tasks and is overall 1.4% be	✓	0.23
CoT prompting provides double-digit improvements for all three models (PaLM, InstructGPT, and Codex) in Table 2.	×	0.05
Codex with CoT prompting outperforms the average human-rater score on 17 out of 23 tasks, compared to 5 out of 23 tasks	✓	0.23
Codex with CoT prompting outperforms the average human-rater by more than 6%, but it still lags behind the best human-ra	✓	0.22
CoT prompting has negative or zero performance gain for text-ada-001 to text-curie-002, but the performance delta betwee	×	0.07
For the PaLM models, CoT prompting has negative performance gain for the smallest model size (8B), but the performance i	×	0.06
CoT is an emergent prompting strategy that requires sufficiently large models.	×	0.05

## References

- <http://arxiv.org/abs/2509.25160v1>
- <http://arxiv.org/abs/2210.09261v1>
- <http://arxiv.org/abs/2502.19187v2>