

# Impact of SFT and DPO on OPT-350M Zero-Shot Reasoning in Big-Bench Hard

Assignee Research

June 12, 2026

## Abstract

BIG-Bench (Srivastava et al., 2022) is a diverse evaluation suite that focuses on tasks believed to be beyond the capabilities of current language models. Language models have already made good progress on this benchmark, with the best model in the BIG-Bench paper outperforming average reported human-rater results on 65% of the BIG-Bench tasks via few-shot prompting. But on what tasks do language models fall short of average human-rater performance, and are those tasks actually unsolvable by current language models? In this work, we focus on a suite of 23 challenging BIG-Bench tasks which we

## 1 Introduction

This paper examines: Challenging BIG-Bench Tasks and Whether Chain-of-Thought Can Solve Them. Research question: To what extent does the combination of SFT and DPO degrade the zero-shot reasoning capabilities of OPT-350M on the Big-Bench Hard suite relative to the base model?.

## 2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 9.2/10.

## 3 Results

11 papers retrieved. 10 claims extracted; 10 independently verified. Quality review score: 9.2/10.

## 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
PaLM, InstructGPT, and Codex models were evaluated on BBH for answer-only and CoT prompting approaches.	✓	0.21
Answer-only prompting typically underestimates language model performance on challenging tasks requiring multiple reason	✓	0.22
In the BIG-Bench paper, none of the evaluated models, including PaLM 540B, outperformed human-rater baselines on any of	✓	0.31
PaLM 540B with answer-only prompting outperforms the average human-rater on 6 out of 23 BBH tasks and is overall 1.4% be	✓	0.36
CoT prompting provides double-digit improvements for all three models (PaLM, InstructGPT, and Codex) in Table 2.	✓	0.22
Codex with CoT prompting outperforms the average human-rater score on 17 out of 23 tasks, compared to 5 out of 23 tasks	✓	0.33
Codex with CoT prompting outperforms the average human-rater by more than 6%, but it still lags behind the best human-ra	✓	0.31
CoT prompting has negative or zero performance gain for text-ada-001 to text-curie-002, but the performance delta betwee	✓	0.42
For the PaLM models, CoT prompting has negative performance gain for the smallest model size (8B), but the performance i	✓	0.40
CoT is an emergent prompting strategy that requires sufficiently large models.	✓	0.15

## References

- <http://arxiv.org/abs/2210.09261v1>
- <http://arxiv.org/abs/2510.01616v1>
- <http://arxiv.org/abs/2305.14947v2>