

Invalid Logic, Equivalent Gains: The Bizarreness of Reasoning in Language Model Prompting

Rylan Schaeffer*¹ Kateryna Pistunova*² Samar Khanna*¹ Sarthak Consul* Sanmi Koyejo¹

¹Computer Science, Stanford ²Physics, Stanford



Summary

- Language models can be prompted to reason through problems in a manner that improves performance
- Why* such prompting improves performance is unclear
- Wang et al. [4] recently showed that *logically invalid* Chain-of-Thought (CoT) prompting [1, 5] improves performance almost as much as logically valid CoT prompting

Andrew Lampinen
@AndrewLampinen

Would be interesting to see if these results are different on more unusual tasks like the BIG Bench ones used in arxiv.org/abs/2210.09261 (by @jasonwei et al.)—in keeping with the interpretation we gave in our explanations work (arxiv.org/abs/2204.02329).

Boshi Wang @BoshiWang2 · Dec 20, 2022

Chain-of-Thought (CoT) prompting improves large LMs on complex reasoning. But what aspects of the CoT prompt matter? Are valid reasoning steps important?

Check out our paper: arxiv.org/abs/2212.10001! Work w/ @sewon_min, @xiangdeng1, @mickeysjm, You Wu, @LukeZettlemoyer, @hhsun1

Jason Wei @jasonwei · Dec 21, 2022

Oh yeah. There is no way the navigate tasks works if the explanations are wrong i feel like...

- Critics responded Wang et al.'s finding was based on too few & too easily solved tasks to draw conclusions
- To resolve this dispute, we test whether *logically invalid* CoT prompts offer the same performance gains on the hardest tasks in BIG-Bench [2], termed BIG-Bench Hard (BBH) [3]
- Logically invalid CoT prompts **DO** achieve similar performance gains on BBH
- We also discover some CoT prompts used by previous works contain logical errors

Background: BIG-Bench [2] & BIG-Bench Hard (BBH) [3]

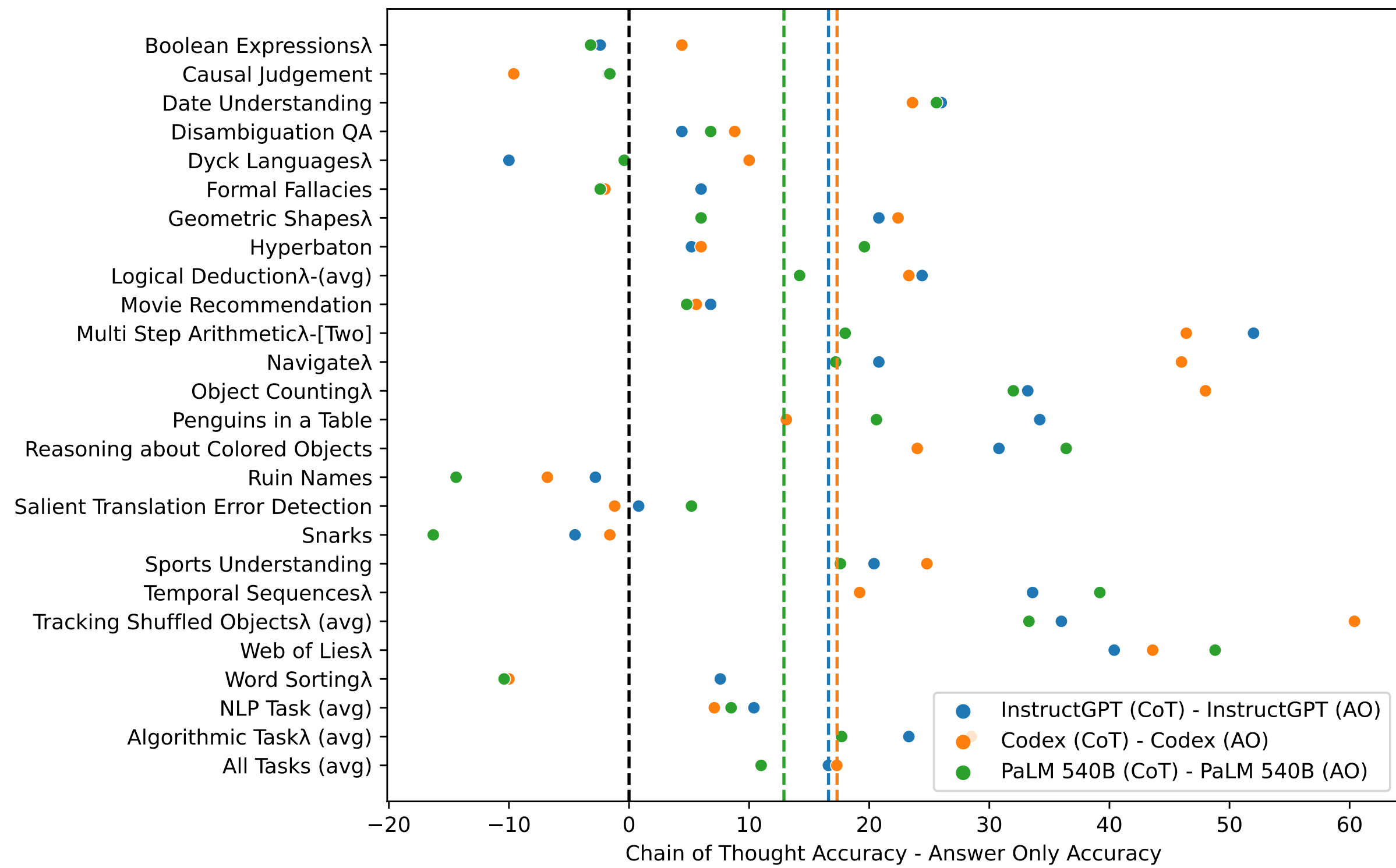
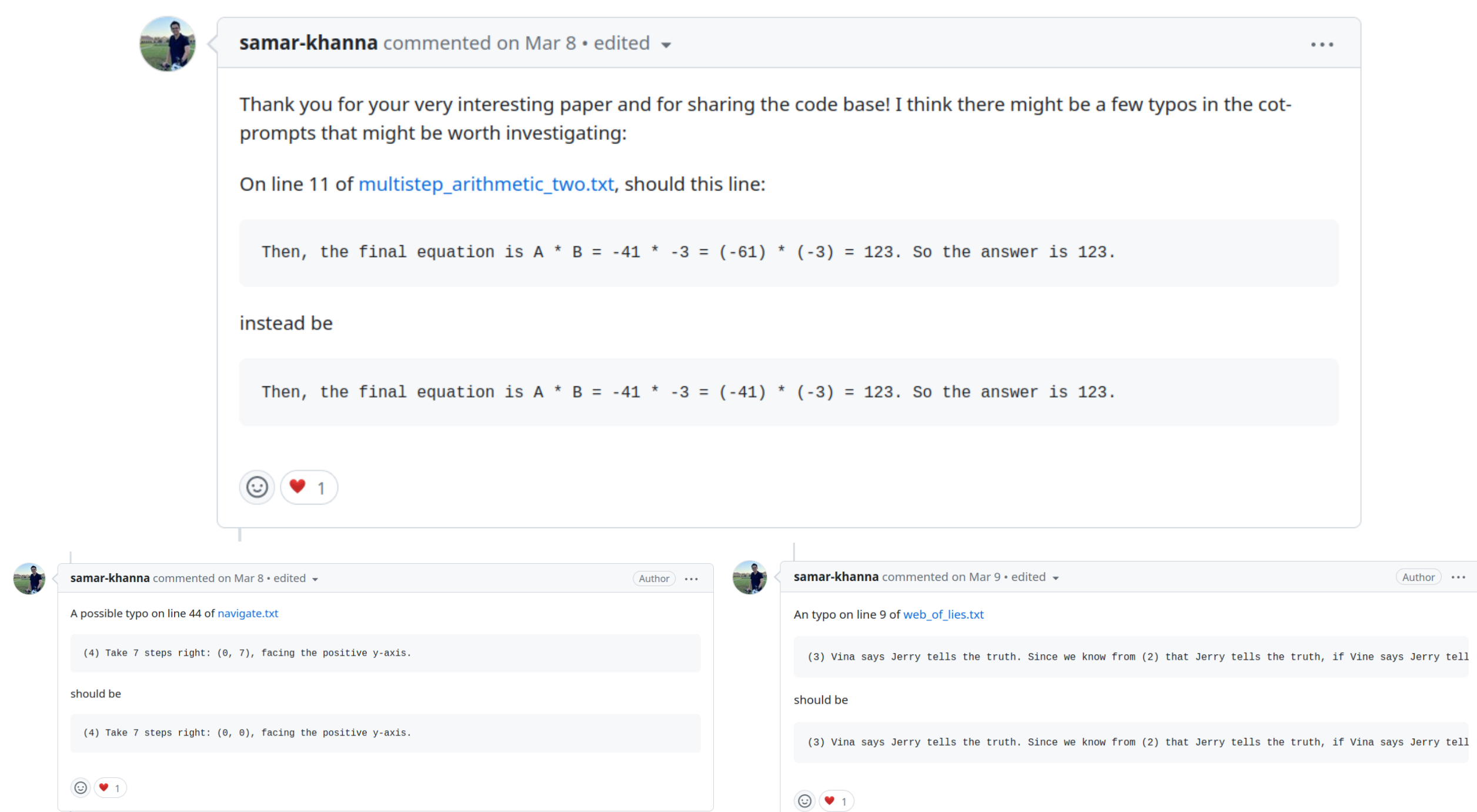
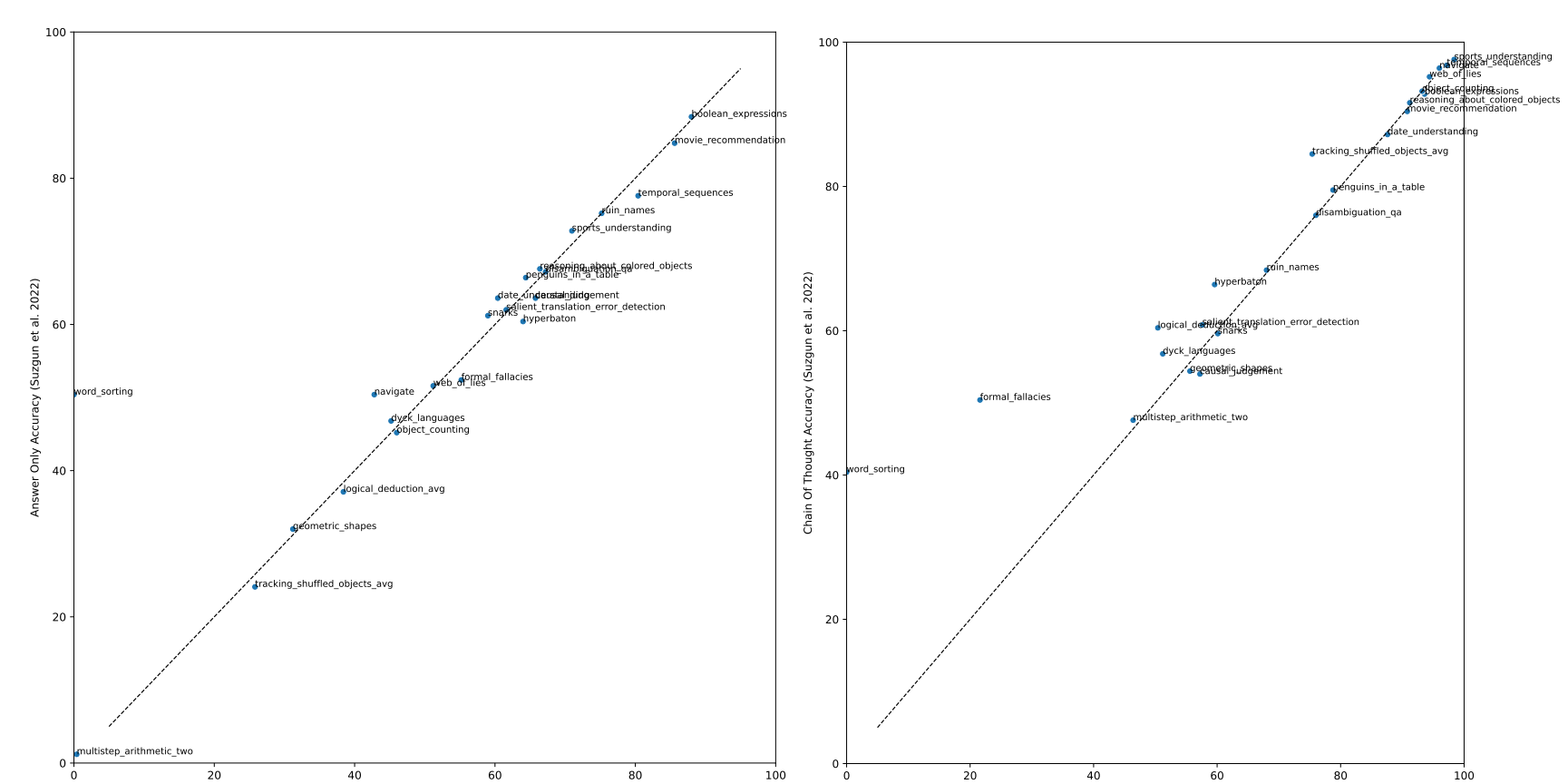


Figure 1. Chain-of-Thought (CoT) prompting significantly outperforms answer-only (AO) prompting on BIG-Bench Hard [3] = 23 of the hardest tasks in Beyond the Imitation Game Benchmark [2].

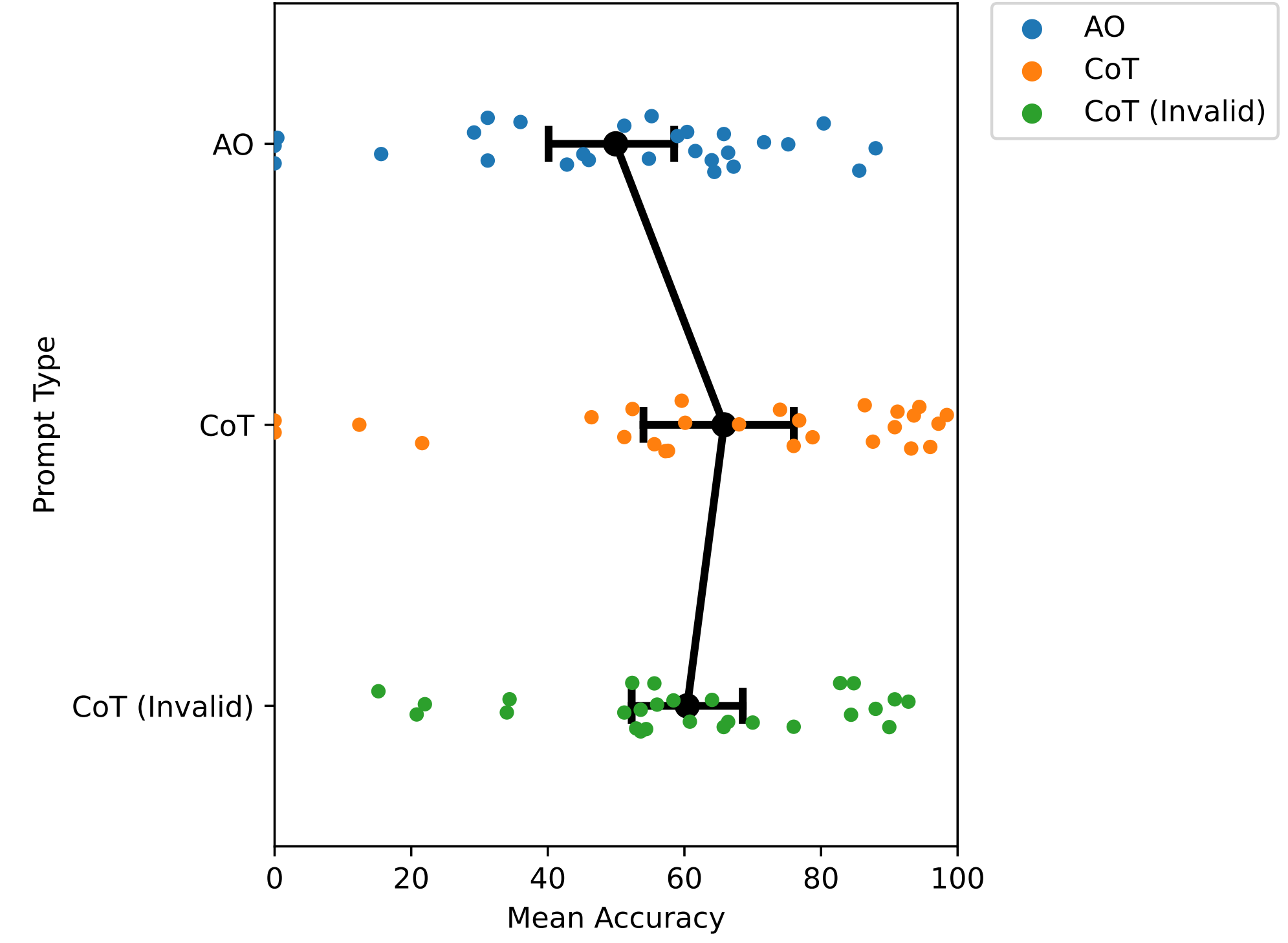
Discovery: BIG-Bench Hard's CoT Prompts Contain Errors!



Reproducing BIG-Bench Hard's Results



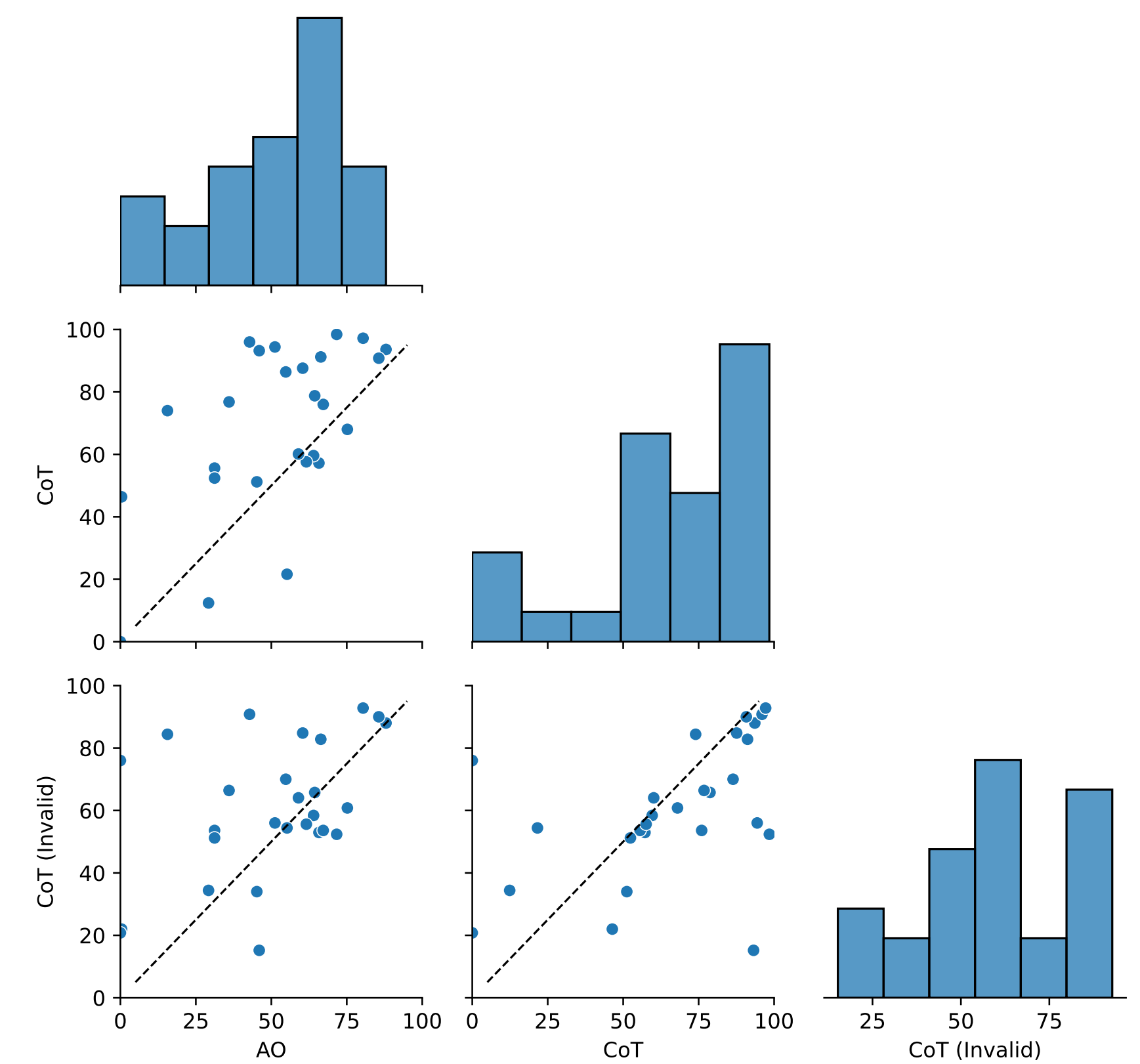
Logically Invalid CoT Prompting Almost Matches CoT on BBH



Note: We do *not* optimize logically invalid CoT prompts whatsoever. Write, run, done.

Example: Evaluate the result of a random Boolean expression. Q: not ((not not True)) is

Answer Only (AO)	False
Scratchpad/Chain-of-Thought (CoT)	A: Let's think step by step. Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of operations from highest priority to lowest priority is "not", "and", "or", respectively. We first simplify this expression "Z" as follows: "Z = not ((not not True)) = not ((A))" where "A = not not True". Let's evaluate A: A = not not True = not (not True) = not False = True. Plugging in A, we get: Z = not ((A)) = not ((True)) = not True = False. So the answer is False.
Logically-Invalid Chain-of-Thought (Invalid CoT)	A: Let's think step by step. Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of operations from highest priority to lowest priority is "not", "and", "or", respectively. Because English does not permit multiple negatives, the expression "(not not True)" evaluates to "(not True)". The expression "not ((not not True))" therefore evaluates to "not ((not True))". By the same logic, the expression "not ((not True))" simplifies to "not True". In Boolean logic, "not True" is False. So the answer is False.



task	AO	CoT	CoT (Invalid)
boolean_expressions	88.0	93.6	88.0
causal_judgement	65.8	57.2	52.9
date_understanding	60.4	87.6	84.8
disambiguation_qa	67.2	76.0	53.6
dyck_languages	45.2	51.2	34.0
formal_fallacies	55.2	21.6	54.4
geometric_shapes	31.2	55.6	53.6
hyperbaton	64.0	59.6	58.4
logical_deduction_five_objects	31.2	52.4	51.2
logical_deduction_seven_objects	29.2	12.4	34.4
logical_deduction_three_objects	54.8	86.4	70.0
movie_recommendation	85.6	90.8	90.0
multistep_arithmetic_two	0.4	46.4	22.0
navigate	42.8	96.0	90.8
object_counting	46.0	93.2	15.2
penguins_in_a_table	64.4	78.8	65.8
reasoning_about_colored_objects	66.4	91.2	82.8
ruin_names	75.2	68.0	60.8
salient_translation_error_detection	61.6	57.6	55.6
snarks	59.0	60.1	64.0
sports_understanding	71.7	98.4	52.4
temporal_sequences	80.4	97.2	92.8
tracking_shuffled_objects_five_objects	15.6	74.0	84.4
tracking_shuffled_objects_seven_objects	0.0	0.0	76.0
tracking_shuffled_objects_three_objects	36.0	76.8	66.4
web_of_lies	51.2	94.4	56.0
word_sorting	0.0	0.0	20.8

References

- M. Nye, A. J. Andreassen, G. Gur-Ari, H. Michalewski, J. Austin, D. Bieber, D. Dohan, A. Lewkowycz, M. Bosma, D. Luan, et al. Show your work: Scratchpads for intermediate computation with language models, 2021.
- A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shob, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models, 2022.
- M. Suzgun, N. Scales, N. Schärli, S. Gehrmann, Y. Tay, H. W. Chung, A. Chowdhery, Q. V. Le, E. H. Chi, D. Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them, 2022.
- B. Wang, S. Min, X. Deng, J. Shen, Y. Wu, L. Zettlemoyer, and H. Sun. Towards understanding chain-of-thought prompting: An empirical study of what matters, 2022.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models, 2022.