

---

# Can Large Language Models Truly Follow your Instructions?

---

**Joel Jang\***  
KAIST  
joeljang@kaist.ac.kr

**Seongheon Ye\***  
KAIST  
seongheon.ye@kaist.ac.kr

**Minjoon Seo**  
KAIST  
minjoon@kaist.ac.kr

## Abstract

In this work, to test the capabilities of large language models on truly following the given instructions, we evaluate 9 common NLP benchmarks with negated instructions on (1) pretrained LMs (OPT & GPT-3) of varying sizes (125M - 175B), (2) LMs further pretrained to generalize to novel instructions (InstructGPT), (3) LMs provided with few-shot examples, and (4) LMs fine-tuned specifically on negated instructions; all LM types perform worse on negated instructions as they scale and show a huge performance gap between the human performance when comparing the average score on both original and negated instructions. By highlighting a critical limitation of existing LMs and methods, we urge the community to develop new approaches to developing LMs that actually follow the given instructions in order to prevent catastrophic consequences that may occur if we prematurely endow LMs with real-world responsibilities.

## 1 Introduction

Large Language Models (LMs) pretrained on a vast amounts of corpora have shown surprising, even emergent, capabilities of solving various downstream tasks through prompts (instructions) [5, 17, 6, 26, 24]. Previous work has specifically shown LMs can perform *unseen* tasks through multitask fine-tuning on various downstream tasks with prompts [20, 23, 21, 15]. A 540B LM [6] has even shown the capability to act as the “brain” for actual robots, helping them perform different tasks in the real-world [1, 10]. As LMs become more capable of performing real-world tasks and are endowed with responsibilities that may result in real-world consequences, it is more-so important to ensure that LMs actually do what they are instructed to do.

In this work, we test the capabilities of Language Models (LMs) on truly following the given instructions (prompts) by conducting a case study with *negated* instructions; that is, telling the LM NOT to do something as shown by an example in Figure 1. Prior work [8, 22] has shown that LMs (as well as other large pretrained models in different modalities such as DALLE-2 [18]) have a hard time understanding negated prompts and perform the task as if provided with the original prompt. For example, if we prompt DALLE-2 this prompt: “Do not generate a monkey holding a banana”, it will generate an image with a monkey holding a banana. Another speculative example is what if LMs are endowed with the decision making of robots [1, 10] and are instructed the following statement: “Whatever you do, do not kill this person”. If the LM does not understand the concept of *negation*, it will lead to catastrophic consequences.

---

\*denotes equal contribution

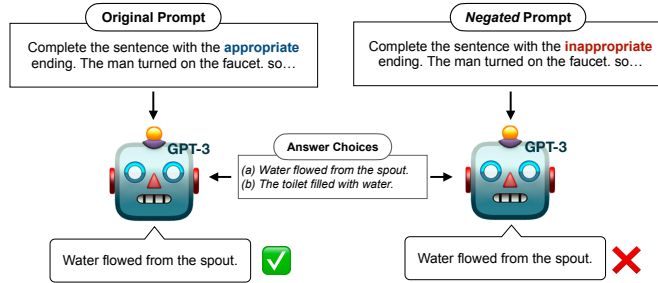


Figure 1: Example of evaluating a sample of the COPA dataset with both the original prompt and the *negated* prompt on GPT-3.

We aim to answer four main questions in this work. (1) How does scaling the size of LMs affect their abilities to understand the concept of negation? (2) Are LMs explicitly trained to follow instructions (InstructGPT) better at understanding negated instructions? (3) Can In-Context Learning (ICL) or Fine-tuning (FT) help mitigate this problem? (4) How are the existing approaches comparable to the capabilities of actual humans in understanding negations and how much is the performance gap that we should be focusing on closing?

The answers to the questions above can be summarized as follows:

- Results show scaling to be inefficient at helping LMs understand *negation*. On the contrary, LMs perform worse as they scale.
- LMs specifically adapted to generalize to novel instructions still suffer from understanding *negation*, despite showing a bit of an improvement.
- ICL helps LMs understand *negation* in only specific scenarios, while FT seems to help in all scenarios. However, FT results in degradation of the original task performance, resulting in a zero-sum game.
- Comparing the existing approaches with the human performance measured by asking 13-year-old humans to do the same task given both the original and *negated* prompts, we show that there is a huge ( $\sim 31.3\%$ ) performance gap to close.

Through this work, we aim to highlight a critical shortcoming of large LMs that should be carefully considered before empowering them with responsibilities that might result in catastrophic real-world consequences, even existential risk <sup>2</sup>.

## 2 Task Description

In this work, we set up a task of evaluating existing LMs capabilities to understand *negated* prompts. We provide the details of how we construct the task in Section 2.1 and the baseline models and methodologies in Section 2.2.

### 2.1 Task Construction

We first choose 9 different datasets categorized into three task types: 3 commonsense reasoning datasets (PIQA [4], ARC-Easy [7], COPA [9]), 3 sentence completion datasets (HellaSwag [25], StoryCloze [14], Lambada [16]), and 3 question answering datasets (WQ [3], NQ [12], TriviaQA [11]). We use the Promptsources Library [2] to find prompts for all of the 9 datasets that show good performance when used to perform tasks with the OPT LMs [26]. Next, we manually *negate* the original prompts. The full list of original and negated prompts are provided in Appendix 5.

<sup>2</sup>As a speculative note, if LLMs do indeed become the foundation of future intelligent agents, it will be critical that LLMs understand the concept of negation. In certain situations, Humans are advised NOT to do something since it is a stronger prevention mechanism (e.g. DO NOT ... signs). LLMs that do not understand the concept of negation and are responsible for making real-world decisions will bring catastrophic results that will be the opposite of human-aligned values.

For evaluation, we sample 300 data instances from each dataset due to the high cost of performing inference with OpenAI API <sup>3</sup>. We use the 300 instances to evaluate *both* the original and negated prompts (a total of 600 data instance inferences for each task). For multi-choice tasks with more than two options such as ARC-Easy, Lambada, and HellaSwag, we consider multiple options to be correct for the negated prompts. For setting up the multiple choice candidates for Lambada, we sampled the other options by choosing a random word from the given input instance. For the 3 QA tasks, we chose the other option by sampling from the answer candidate list from the training set that did not overlap with any of the answer candidates from the test set. We provide the final data instances used for the evaluation via csv files at this link.

## 2.2 Baselines

**Baseline Models** For the main experiments, we use the OPT LMs [26] (125M, 350M, 1.3B, 2.7B, 6.7B, 13B, 30B, 66B, 175B) and GPT-3 [5] LMs (Ada, Babbage, Curie, Davinci) to observe the effect of *scale* on the capabilities of LMs to better follow the given prompts.

**Existing Methodologies** We explore how much existing methodologies can help LMs understand negation by performing the experiments with LMs further adapted to follow instruction (T0 [20] 3B and 11B); InstructGPT [15] Ada, Babbage, Curie, Davinci), with In-Context Learning (ICL) on OPT 66B model with  $k = 2, 4, 8$  shots, and Fine-tuning (FT) with OPT 125M, 350M, and 1.3B. For fine-tuning tasks with  $<10k$  training instances, we train on all of the available training data (with negated prompts) for 5 epochs. For tasks with training data  $>10k$ , we limit the training instance number to 10k and train for 5 epochs as well. We use a fixed learning rate of  $1e-3$ . We use the last model checkpoint for evaluating our task.

**Human Evaluation** We also provide human evaluations on 3 different tasks, one for each task category (COPA for commonsense reasoning, Lambada for sentence completion, and NQ for question answering). We sample 100 out of the 300 instances: 50 instances with the original prompt and 50 instances with the negated prompt (non-overlapping) and evaluate them on three 13-year-old humans. We did this to quantify exactly how much these LLMs perform poorly on understanding the concept of negation compared to humans that we hypothesized could easily understand and perform negated prompts, even 13-year-old humans <sup>4</sup>.

## 3 Experimental Results

### 3.1 The Effect of Scale

We show the results of evaluating all scales of GPT-3 LMs on our task setup in Figure 2. Since the OPT LMs show the same trend, we show the results in Appendix 6. We find that for all tasks (commonsense reasoning, sentence completion, question answering), an *inverse* scaling law is shown: larger LMs tend to perform worse on negated prompts. This result is very unexpected considering that the zero-shot performance of LMs improves as the size of the LMs increase as shown via the original prompt performance [5, 24]. This leads to a flat line performance for the **average** of negative and positive prompts that is  $\sim 50\%$  for all of the tasks. In other words, this means that the LMs could not find any distinction between the original and the negated prompts, treating them as identical instructions when in reality, those prompts are asking the LM to do opposite things.

**A Conjecture on Why Inverse Scaling Exists** We conjecture that this *inverse* scaling law of negated prompts is caused by a bias from the pretraining corpora towards favoring the original prompts to the negated prompts. Therefore, unless we control the balance between positive and negative texts from the pretraining corpora, which is practically infeasible, this problem will be difficult to solve. We believe it will be especially more difficult for large LMs, since they are powerful language modeling representers, meaning that it would be harder to make the LM *revert* the label prediction by focusing on a single negative word "not" or other negation words that are rarely observed in the training corpora. Large LMs would treat the negation word as a grammatical error and perform language modeling of the positive texts. This is based on the assumption that large

<sup>3</sup><https://openai.com/api/>

<sup>4</sup>We received advice and guidance from the National Human Rights Commission of Korea for setting up the task evaluation on the minors and received parental consent.

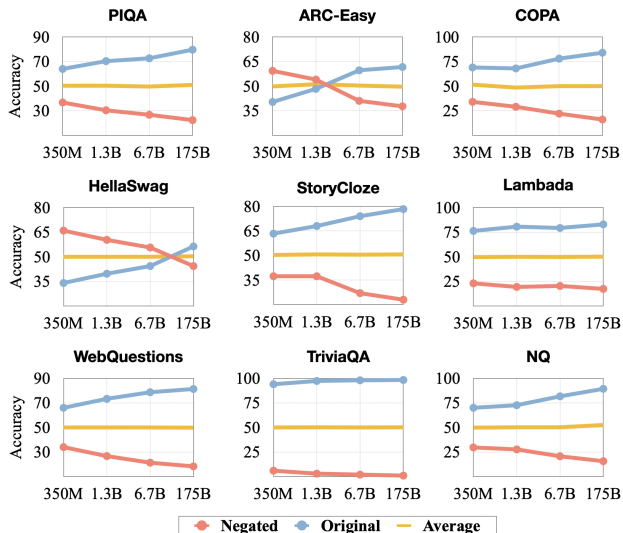


Figure 2: Zero-shot task performance of 9 datasets on GPT-3 across different model scales (350M, 1.3B, 6.7B, 175B).

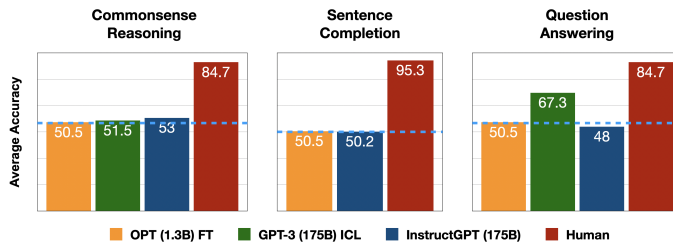


Figure 3: Results from each of the 3 task-type show how the existing methods compare with the human performance on the average score of both original and negated prompts. We do not report GPT-3 (175B) ICL results for Sentence Completion (Lambada) because  $k=8$  shots exceeded the token limit of the OpenAI API.

pretrained LMs might be just probabilistic models instead of actually *actively* learning the linguistic knowledge, which is also questioned in previous works [13, 19] as well. However, thorough analysis and experiments are further needed in order to validate this conjecture, which we leave for future work.

### 3.2 Existing Methodologies

We show a subset of the experiments performed with other existing methodologies and the performance gap compared to the human evaluation in Figure 3. As shown in Figure 3, existing methods, except for GPT-3 (175B) ICL for question answering, show  $\sim 50\%$  performance on the average of original and negated prompts. Considering the results of zero-shot GPT-3 shown in Section 3.1 and the results shown in this section where the best performing method still has an average performance gap of 31.3% compared with the human performance, the current limitations of existing LMs on precisely understanding and following the given prompts clearly exist. We provide the full experimental results on *negated* prompts of T0 (3B), T0 (11B), all scales of InstructGPT, ICL OPT 66B with  $k = 2, 4, 8$  and FT of OPT 125M, 350M, and 1.3B in Appendix 7.

## 4 Closing

Large LMs have taken the research community on an expeditious journey in the last 2 years, achieving past average human performance on many NLP benchmarks [6], and even extending its use-case to act as the *brain* for actual robots [1, 10]. As the real-world use cases of large LMs are widened, the research community should carefully consider if the LMs precisely understand the given instructions or if they are highly biased on the distribution of the pretraining corpora. We close our case study by

urging the community to develop new methodologies for creating truly instruction-following LMs before relying on their capabilities for making real-world decisions.

## Acknowledgements

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grants funded by the Korea government (MSIT) (No.2022-0-00113, Developing a Sustainable Collaborative Multi-modal Lifelong Learning Framework, 80%; No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST), 10%; No.2021-0-02068, Artificial Intelligence Innovation Hub, 10%).

## References

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [2] Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-david, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Fries, Maged Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Dragomir Radev, Mike Tian-jian Jiang, and Alexander Rush. PromptSource: An integrated development environment and repository for natural language prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 93–104, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [3] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- [4] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [7] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [8] Allyson Ettinger. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48, 2020.
- [9] Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 394–398, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics.
- [10] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.

- [11] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics.
- [12] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019.
- [13] Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- [14] Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California, June 2016. Association for Computational Linguistics.
- [15] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [16] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [17] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- [18] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [19] Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. Impact of pretraining term frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*, 2022.
- [20] Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2021.
- [21] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. Benchmarking generalization via in-context instructions on 1,600+ language tasks. *arXiv preprint arXiv:2204.07705*, 2022.
- [22] Albert Webson and Ellie Pavlick. Do prompt-based models really understand the meaning of their prompts? *arXiv preprint arXiv:2109.01247*, 2021.
- [23] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- [24] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.

Table 1: Full list of original and negated prompts for the 9 evaluation datasets.

Dataset	Type	Prompt
PIQA	Original	Generate the correct solution to accomplish the following goal {goal} {option}
	Negated	Generate the incorrect solution to accomplish the following goal: {goal} {option}
ARC-Easy	Original	Generate the correct answer to the following question. Question: {question} Answer is {option}
	Negated	Generate the incorrect answer to the following question. Question: {question} Answer is {option}
COPA	Original	Complete the sentence with the appropriate ending. {premise} because(so).. {option}
	Negated	Complete the sentence with an inappropriate ending. {premise} because(so).. {option}
HellaSwag	Original	Complete the sentence with an appropriate ending: {input} {option}
	Negated	Complete the sentence with an inappropriate ending: {input} {option}
StoryCloze	Original	Generate a natural ending for the following story: {4 sentences}. option
	Negated	Generate an unnatural ending for the following story: {4 sentences}. {option}
Lambada	Original	Please generate a natural ending following the given chunk of text. {text} {option}
	Negated	Please generate an unnatural ending following the given chunk of text. {text} {option}
WQ	Original	Give me a possible correct answer to the question {question}. {option}
	Negated	Give me a possible incorrect answer to the question {question}. {option}
TriviaQA	Original	What is a correct answer to the following question? Question: {question} Answer: {option}
	Negated	What is an incorrect answer to the following question? Question: {question} Answer: {option}
NQ	Original	The goal is to predict a correct English answer string for an input English question. Question : {question} Answer: {option}
	Negated	The goal is to predict an incorrect English answer string for an input English question. Question : {question} Answer: {option}

- [25] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- [26] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

## 5 Full List of Prompts

Table 1 shows the full list of original and negated instructions used on all of the 9 datasets.

## 6 OPT Results

We show the main results of all scales of the OPT LMs in Figure 4. Same with the results shown on the GPT-3 LMs, performance worsens as the LMs scale larger for the *negated* prompts, resulting in an  $\sim 50\%$  for the average score.

## 7 Evaluation of Different LMs on *Negated* Prompts

Figure 5-13 shows the performance of the existing methods on the *negated* prompts. First thing to note is that for all of the tasks, instruction following LMs (T0 and InstructGPT) also show an inverse scaling law where the larger LMs perform worse. Furthermore, while FT seems to help LMs understand *negation*, as shown in Figure 3, the average score of both original and negated prompts is  $\sim 50\%$ , which means that the mitigation was a result of a trade-off of performance degradation on the original prompts, resulting in a zero-sum game.

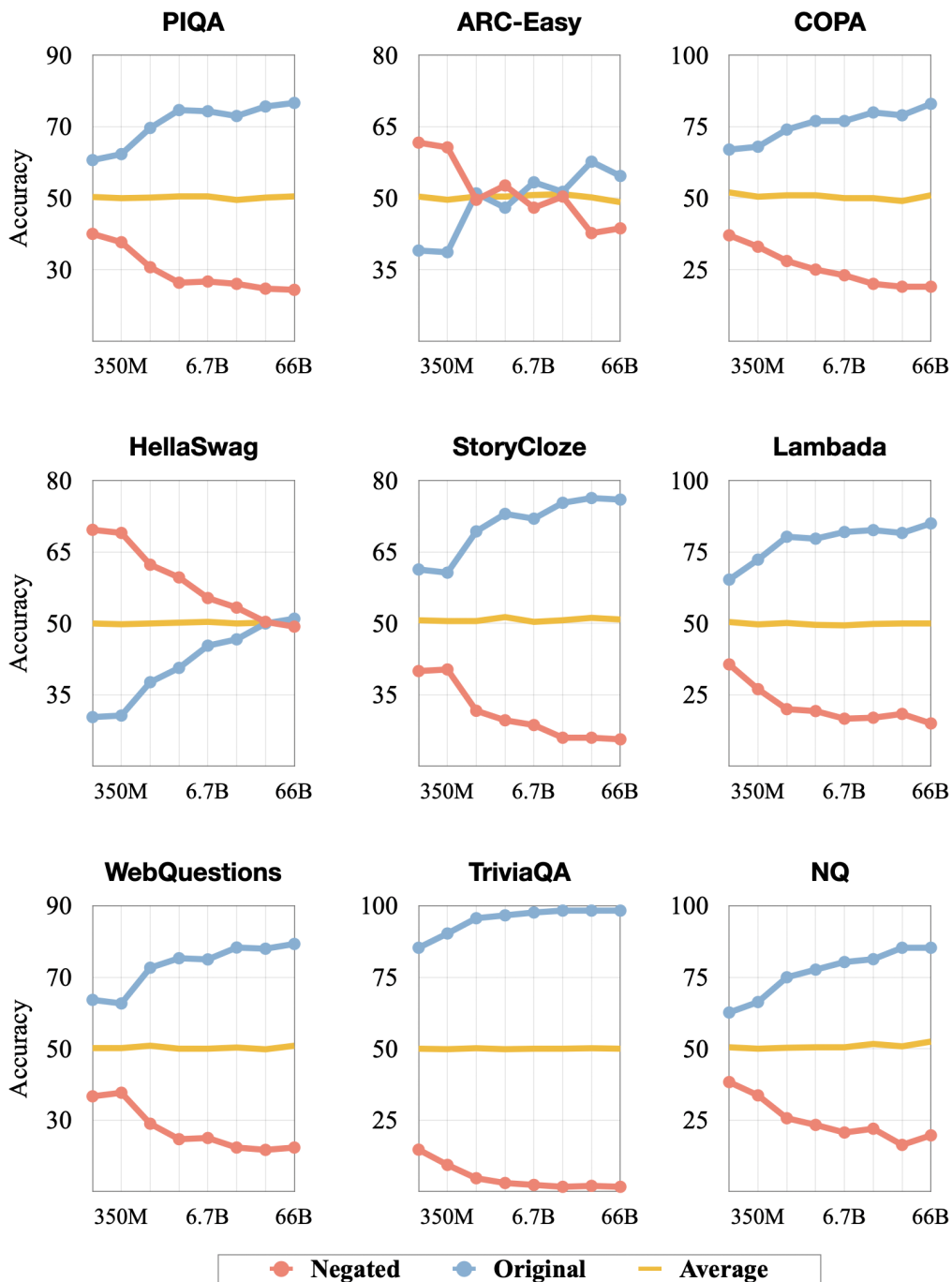


Figure 4: Zero-shot task performance of 9 datasets of OPT across different model scales (125m, 350m, 1.3b, 2.7b, 6.7b, 13b, 30b, 66b).



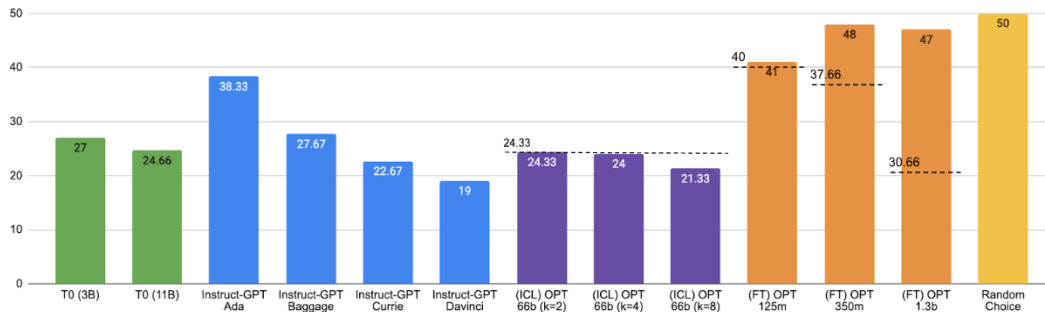


Figure 5: PIQA

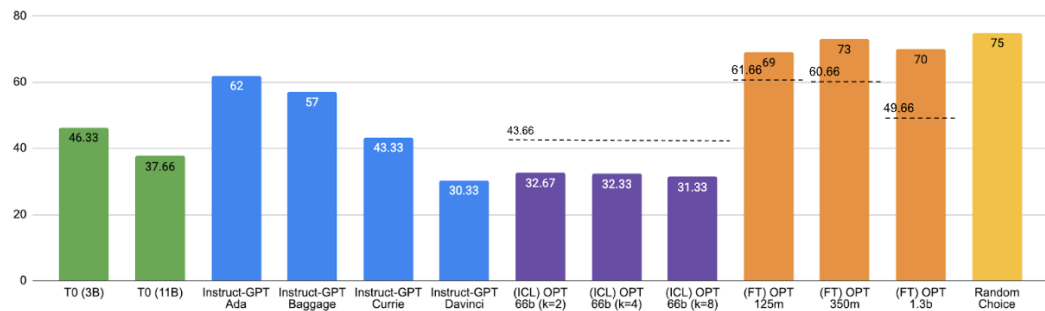


Figure 6: ARC-Easy

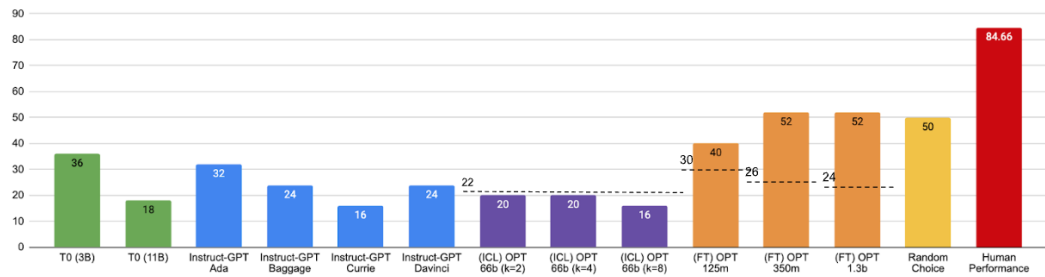


Figure 7: COPA

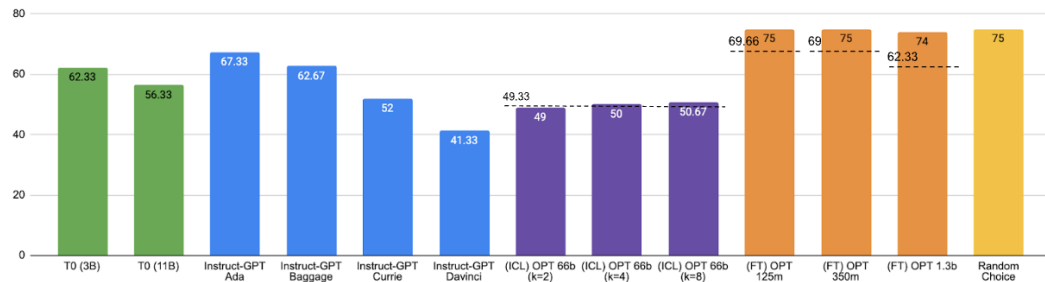


Figure 8: HellaSwag

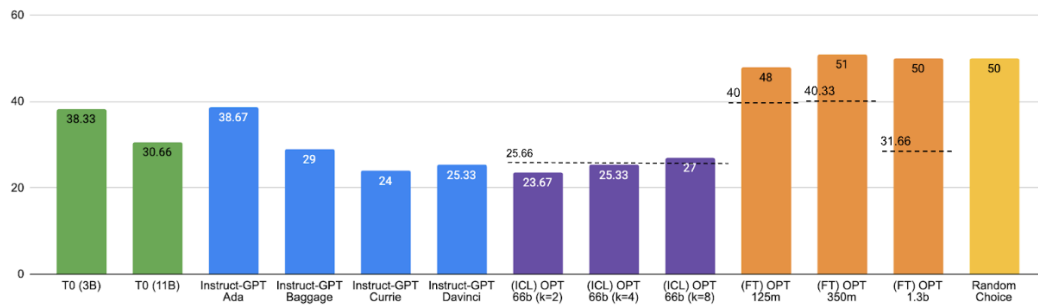


Figure 9: StoryCloze

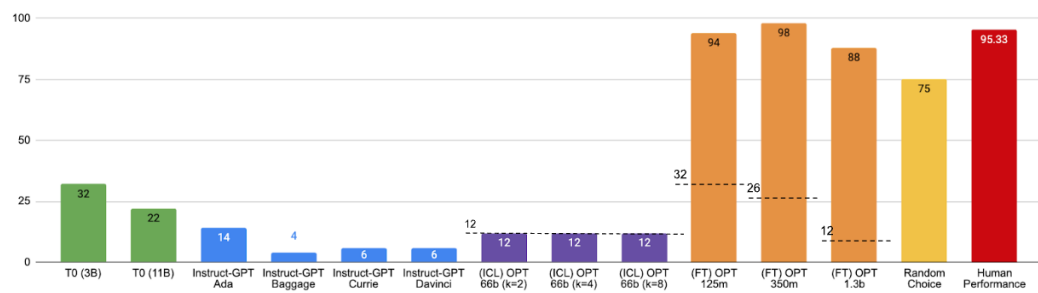


Figure 10: Lambada

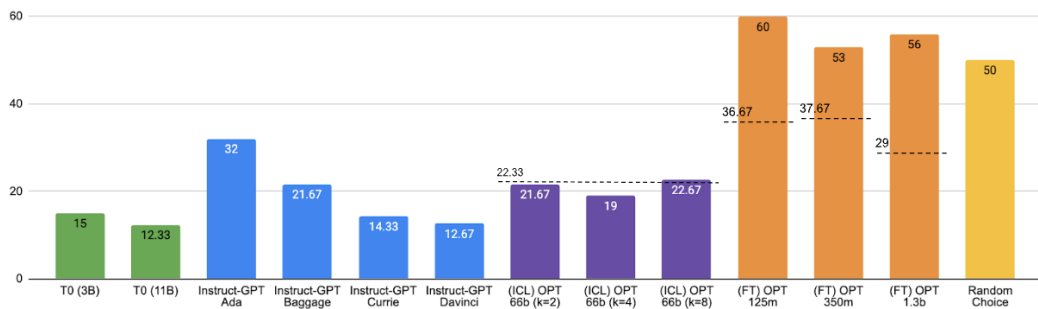


Figure 11: Web Questions

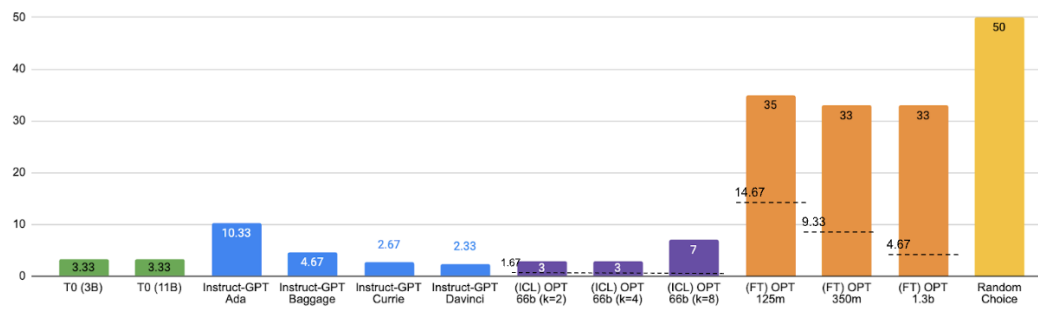


Figure 12: TriviaQA

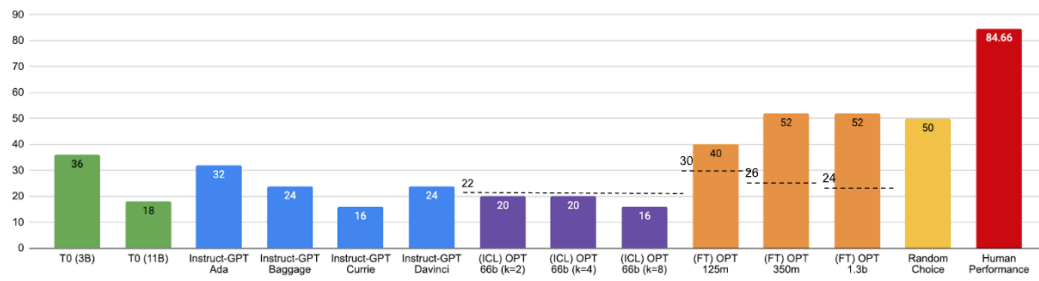


Figure 13: Natural Questions