

# Chain-of-Thought Prompting Elicits Reasoning in Large Language Models



**Paper link:** <https://arxiv.org/pdf/2305.15507.pdf>



**Related Article:** [https://open.substack.com/pub/aiguide/p/can-large-language-models-reason?r=b0ibt&utm\\_campaign=post&utm\\_medium=web](https://open.substack.com/pub/aiguide/p/can-large-language-models-reason?r=b0ibt&utm_campaign=post&utm_medium=web)

## Brief Summary

Chain of thought as a way to improve the scope of what LLMs can do. What is chain of thought? Adding intermediate steps in the form of natural language to improve the model performance. Does not answer much about reasoning. The paper basically says that the capabilities of LLMs expand in various domains such as arithmetic, commonsense, symbolic reasoning as we increase the size of the model and use chain of thought.

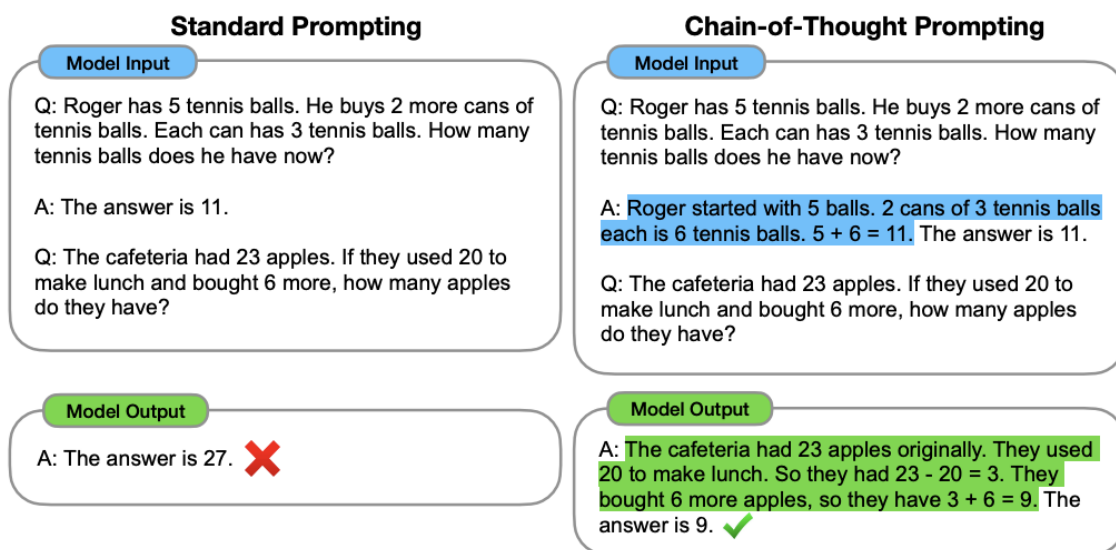


Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

## Detailed summary

Natural Language Processing has been revolutionized by the use of LLMs. If you are using any AI product these days, such as ChatGPT, chances are that LLM is the technology behind it. It has been observed that as we scale up the model, i.e. increase the number of model parameters, increase the amount of training data, the model helps in reasoning but is still not enough for arithmetic tasks. Arithmetic reasoning is a task which is easy for humans but models suffer a lot in them. This paper makes use of natural language to improve performance on arithmetic tasks. The paper takes advantage of few-shot learning, which means that we do not need to fine tune model to a specific task (arithmetic reasoning is this case), but just providing a few sample input-output exemplar demonstrating the task will make the model learn. Few shot learning is a very helpful methodology which saves the cost of creating high set of rational datasets to be used for fine tuning the model. The method proposed by the authors consists of a prompt<input, chain of thought, output>

Chain of thought is an intermediate step that explain the process involved for taking input to the output. It has been shown in the paper that sufficiently large language models can generate chains of thought if demonstrations of chain-of-thought reasoning are provided in the exemplars for few-shot prompting. Chain of thought allows the

model to decompose the problem into multi step so that computation can be allotted to problems that require more reasoning steps. This is also a very interesting method as it explains the model and is applicable to any task that humans can solve using language

The benchmark mathematical problem sets considered for the paper: GSM8K, SVAMP, ASDiv, AQUa, MAWPS

The paper mentions two types of prompting -

Standard prompting: Consider standard few-shot prompting in which language model is given in-context exemplars of input–output pairs before outputting a prediction for a test-time example. There is no chain of thought input here.

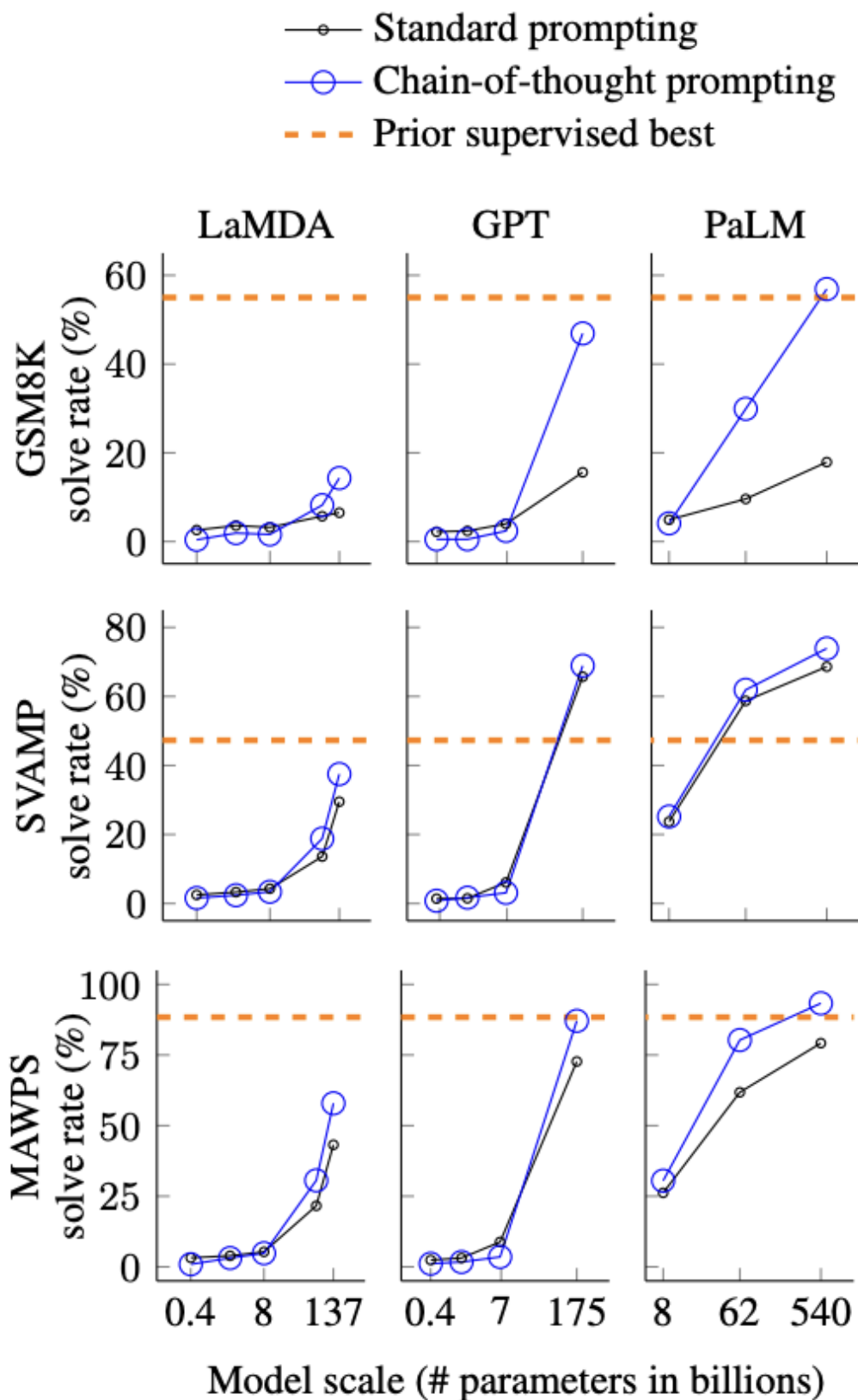
Chain of thought prompting: Augment each exemplar with chain of thought in few shot learning.

<p><b>Math Word Problems (free response)</b></p> <p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. <math>5 + 6 = 11</math>. The answer is 11.</p>	<p><b>Math Word Problems (multiple choice)</b></p> <p>Q: How many keystrokes are needed to type the numbers from 1 to 500? Answer Choices: (a) 1156 (b) 1392 (c) 1480 (d) 1562 (e) 1788</p> <p>A: There are 9 one-digit numbers from 1 to 9. There are 90 two-digit numbers from 10 to 99. There are 401 three-digit numbers from 100 to 500. <math>9 + 90(2) + 401(3) = 1392</math>. The answer is (b).</p>	<p><b>CSQA (commonsense)</b></p> <p>Q: Sammy wanted to go to where the people were. Where might he go? Options: (a) race track (b) populated areas (c) desert (d) apartment (e) roadblock</p> <p>A: The answer must be a place with a lot of people. Race tracks, desert, apartments, and roadblocks don't have a lot of people, but populated areas do. So the answer is (b).</p>
<p><b>StrategyQA</b></p> <p>Q: Yes or no: Would a pear sink in water?</p> <p>A: The density of a pear is about <math>0.6 \text{ g/cm}^3</math>, which is less than water. Thus, a pear would float. So the answer is no.</p>	<p><b>Date Understanding</b></p> <p>Q: The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date 10 days ago in MM/DD/YYYY?</p> <p>A: One day after 06/01/1943 is 06/02/1943, so today is 06/02/1943. 10 days before today is 05/23/1943. So the answer is 05/23/1943.</p>	<p><b>Sports Understanding</b></p> <p>Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."</p> <p>A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.</p>
<p><b>SayCan (Instructing a robot)</b></p> <p>Human: How would you bring me something that isn't a fruit?</p> <p>Explanation: the user wants something to eat that isn't a fruit. An energy bar is not a fruit, so I will bring the user an energy bar.</p> <p>Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().</p>	<p><b>Last Letter Concatenation</b></p> <p>Q: Take the last letters of the words in "Lady Gaga" and concatenate them.</p> <p>A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.</p>	<p><b>Coin Flip (state tracking)</b></p> <p>Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?</p> <p>A: The coin was flipped by Maybelle. So the coin was flipped 1 time, which is an odd number. The coin started heads up, so after an odd number of flips, it will be tails up. So the answer is no.</p>

The evaluation is done on 5 large models namely - GPT-3, LaMDA, UL2 20B, Codex

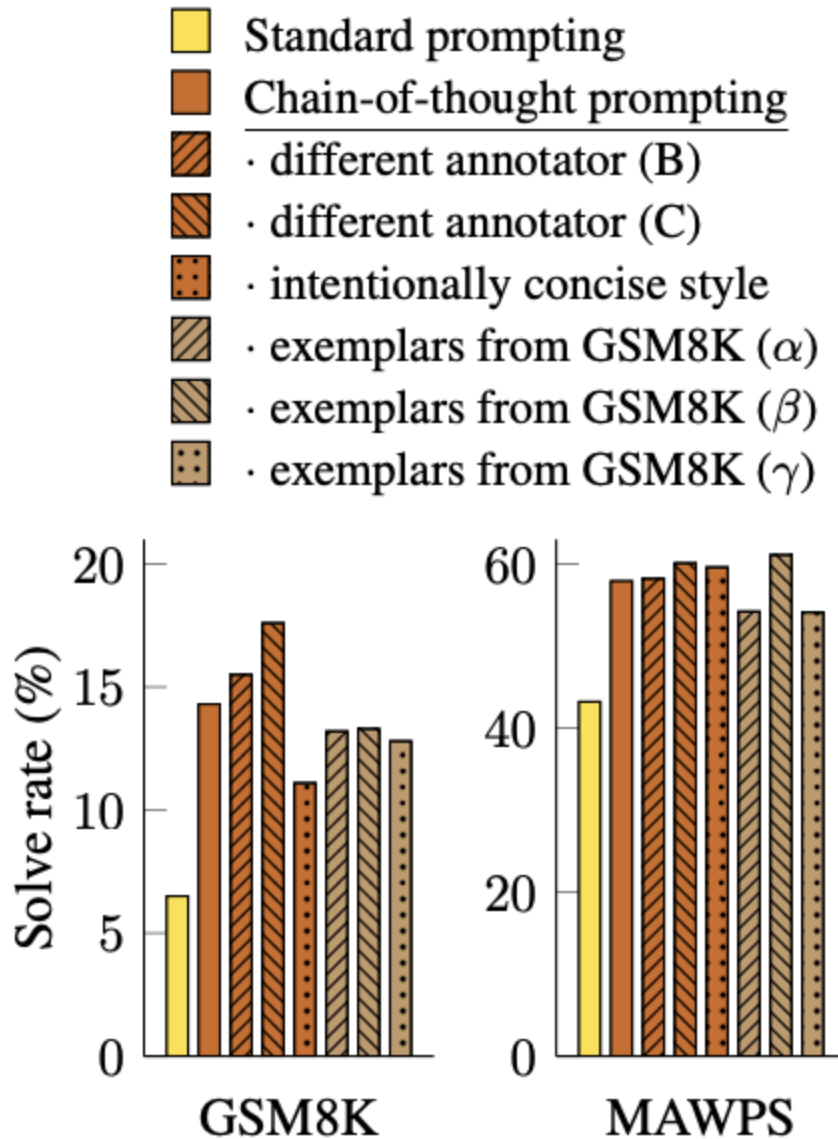
## RESULTS -

1. Not useful for smaller models as the performance is not impacted, qualitatively found that models of smaller scale produced fluent but illogical chains of thought, leading to lower performance than standard prompting
2. Chain-of-thought prompting has larger performance gains for more-complicated problems.
3. Scaling improves the chain of thought i.e. larger the model size, the better quality chain of thoughts produced.



The diagram above clearly shows that for larger model sizes (model scale), chain-of-thought prompting surpasses the benchmark. While standard prompting does improve performance as the size increases, it falls short of both the benchmark and the chain-of-thought prompting method. It can also be observed that when the model scale is less, the performance of chain-of-thought prompting is similar to that of standard prompting because of the quality of chain of thoughts generated by small models.

The paper also mentions different ablation studies that highlight how chain-of-thought is the only variable in this work responsible for improving the model performance as the model scales up. Feel free to refer the paper to understand the ablation study methods used.



The above graph is for LaMDA trained on 137Billion parameters. It can be clearly seen that all the chain-of-thought methods perform better than the standard prompting although they differ by some margin owing to different thought styles written by different persons. This implies that successful use of chain of thought does not depend on a particular linguistic style. In addition to robustness to annotators, independently-written chains of thought, different exemplars, and various language models, the paper also finds that chain-of-thought prompting for arithmetic reasoning is robust to different exemplar orders and varying numbers of exemplars

Apart from mathematic reasoning tasks, the paper also shows the use of chain-of-thought prompting on Commonsense reasoning and Symbolic reasoning and shows that chain-of-thought method does take the model evaluation past the previous benchmark results.