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

Presented by:
Yicong Li, Ziwei Gu, Shivam Raval, Hongwen Song

# Challenges in LLMs

- Scaling up model size alone has not proved sufficient for achieving high performance on challenging tasks, such as arithmetic, commonsense, and symbolic reasoning.
- Large language models still have limitations in their ability to reason and understand the context of a situation.

# Previous work

- Models have the ability to generate neutral language intermediate steps by training from scratch (Ling et al., 2017) or finetuning a pretrained model (Cobbe et al., 2021).
  - costly to create a large set of high quality rationales
- Large language models offer the exciting prospect of in-context few-shot learning via *prompting* (Brown et al., 2020)
  - works poorly on tasks that require reasoning abilities,
  - often does not improve substantially with increasing language model scale

# Contribution

- This work explores the ability of language models to perform few-shot prompting for reasoning tasks, given a prompt that consists of triples: ⟨input, *chain of thought*, output⟩.
    - Chain-of-thought: a series of intermediate natural language reasoning steps that lead to the final output (*chain-of-thought prompting*).
- This work also presents empirical evaluations on arithmetic, commonsense, and symbolic reasoning benchmarks, showing that chain-of-thought prompting outperforms standard prompting.

# Methodology

# How do we solve complicated problems

Decompose the problem into intermediate steps and solve each before giving the final answer

*"After Jane gives 2 flowers to her mom she has 10 . . . then after she gives 3 to her dad she will have 7 . . . so the answer is 7."*

 The paper shows sufficiently large LMs can generate chains of thought if demonstrations of chain-of-thought reasoning are provided as few-shot prompting.

## Standard Prompting

**Model Input**

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?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

## Chain-of-Thought Prompting

**Model Input**

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?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✅

# Benefits

1. Allows models to decompose multi-step problems into intermediate steps so additional computation can be allocated to problems that require more reasoning steps.
2. Interpretable method to understand model's response
3. Potentially applicable to any task that require multi-step approach.
4. Can be induced simply by including examples of chain of thought sequences while prompting.

# Sensitivity to prompt engineering

Chain of thought is robust to:

1. Different annotators writing prompts
2. Annotators without machine learning background writing prompts
3. Different types, ordering and number of prompt examples
4. Different language models (LaMDA, GPT-3, and PaLM)

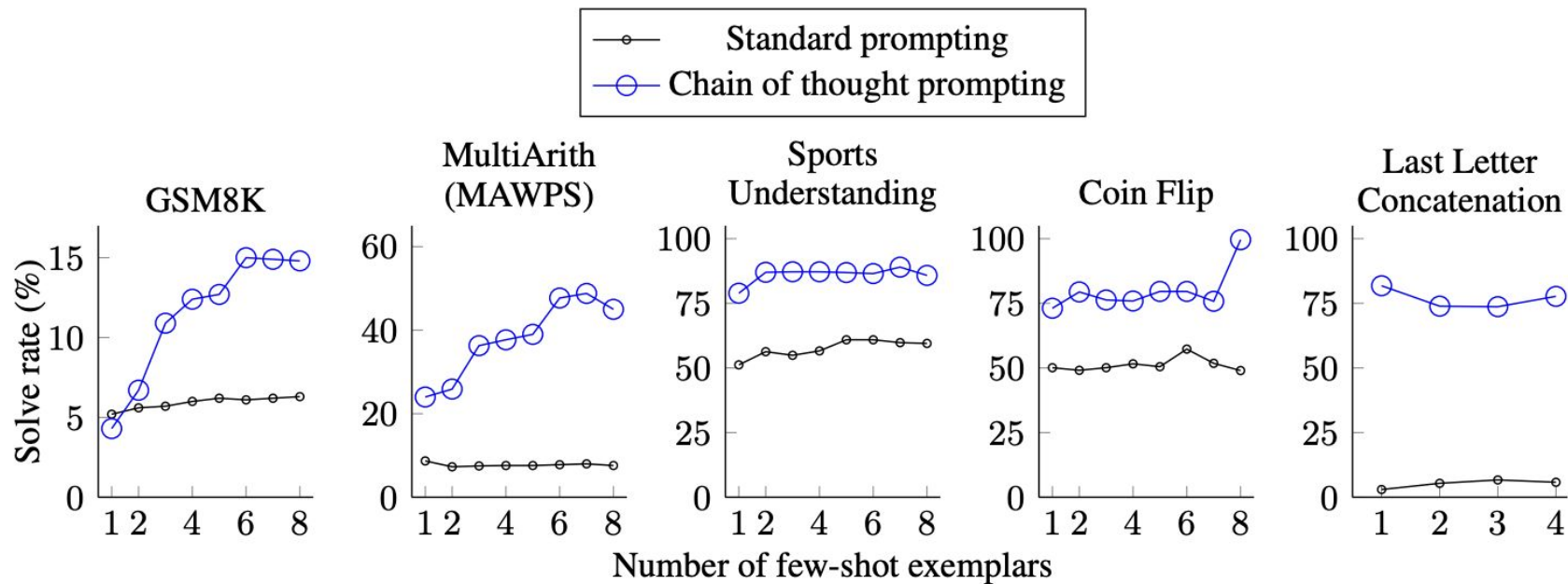However Prompt Engineering improved performance significantly in many cases.

Figure 11: The improvement of chain of thought prompting over standard prompting appears robust to varying the number of few-shot exemplars in the prompt.

Harvard John A. Paulson
School of Engineering
and Applied Sciences

# Experimental Setup

**Benchmarks**: five math word problem benchmarks:
- GSM8K (Cobbe et al., 2021)
- SVAMP (Patel et al., 2021)
- ASDiv (Miao et al., 2020)
- AQuA
- MAWPS (Koncel-Kedziorski et al., 2016)

---

**PROMPT FOR MATH WORD PROBLEMS**

**Q:** There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

**A:** There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6.

Harvard John A. Paulson
**School of Engineering**
and Applied Sciences

# Experimental Setup

**Models**:
- baseline: standard few-shot prompting (Brown et al., 2020)
- chain-of-thought prompting: augment each exemplar in few-shot prompting with a chain of thought for an associated answer, using five different large language models:
    - GPT-3 (Brown et al., 2020)
    - LaMDA (Thoppilan et al., 2022)
    - PaLM
    - UL2 20B (Tay et al., 2022)
    - Codex (Chen et al., 2021)

# Results

**Three key takeaways:**
- chain-of-thought prompting is an emergent ability of **model scale**;
- chain-of-thought prompting has larger performance gains for **more-complicated problems**;
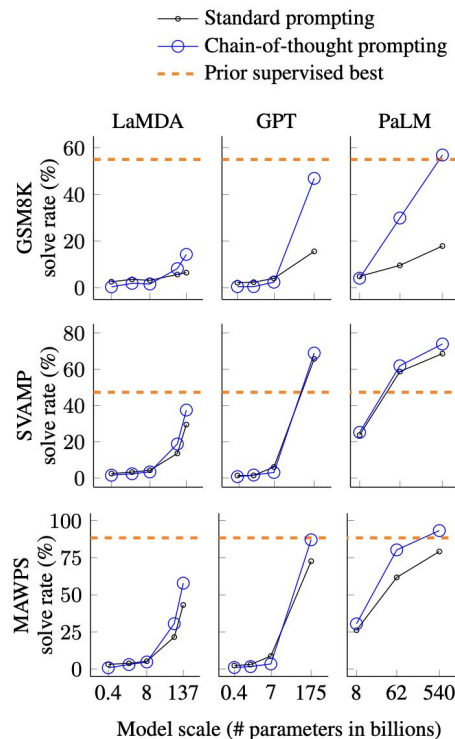- chain-of-thought prompting via GPT-3 175B and PaLM 540B **compares favorably** to prior state of the art.



Figure 4: Chain-of-thought prompting enables large language models to solve challenging math problems. Notably, chain-of-thought reasoning is an emergent ability of increasing model scale.

# Ablation Study

**An ablation study with three variations of chain of thought:**
- Equation only
- Variable compute only
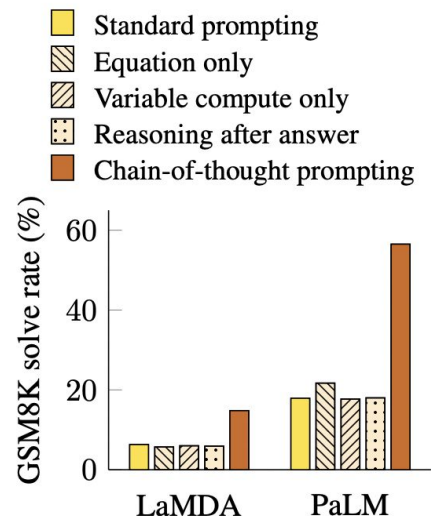- Chain of thought after answer



Figure 5: Ablation study for different variations of prompting using LaMDA 137B and PaLM 540B. Results for other datasets are given in Appendix Table 6 and Table 7.

Harvard John A. Paulsor
**School of Engineering**
and Applied Sciences

# Robustness of Chain of Thought

**Findings:** chain-of-thought prompting for arithmetic reasoning is **robust** to
- annotators
- independently-written chains of thought
- different exemplars
- various language models
- different exemplar orders
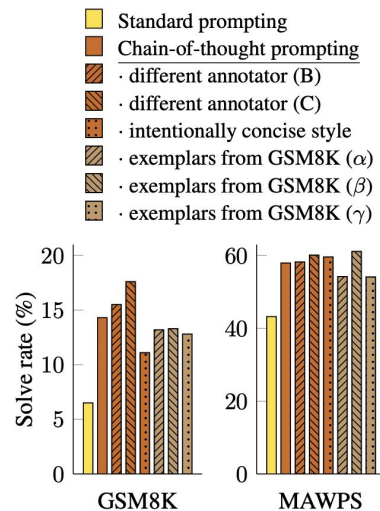- varying numbers of exemplars



Figure 6: Chain-of-thought prompting has variance for different prompt examples (as expected) but outperforms standard prompting for various annotators as well as for different exemplars.

# Commonsense Reasoning

## Math Word Problems (free response)

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?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

## Math Word Problems (multiple choice)

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

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. 9 + 90(2) + 401(3) = 1392. The answer is (b).

## CSQA (commonsense)

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

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).

## StrategyQA

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about 0.6 g/cm^3, which is less than water. Thus, a pear would float. So the answer is no.

## Date Understanding

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?

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.

## Sports Understanding

Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."

A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.

## SayCan (Instructing a robot)

Human: How would you bring me something that isn't a fruit?

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.
Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().

## Last Letter Concatenation

Q: Take the last letters of the words in "Lady Gaga" and concatenate them.

A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.
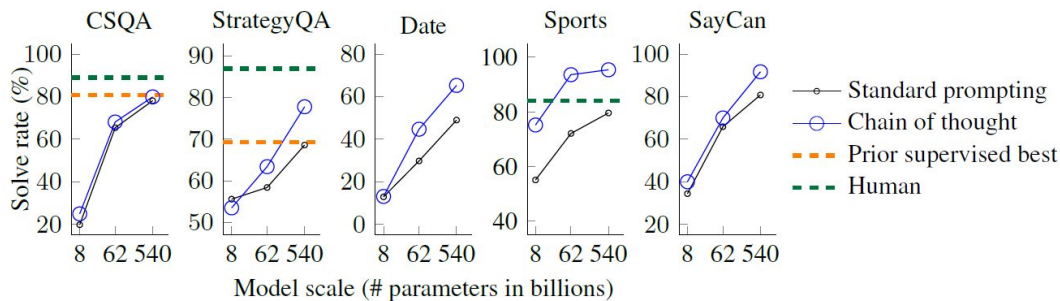
## Coin Flip (state tracking)

Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?

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.

# Results



- Chain of thought prompting works better when the model is sufficiently large.

Table 4: Standard prompting versus chain of thought prompting on five commonsense reasoning benchmarks. Chain of thought prompting is an emergent ability of model scale—it does not positively impact performance until used with a model of sufficient scale.

| Model | | CSQA standard | CoT | StrategyQA standard | CoT | Date standard | CoT | Sports standard | CoT | SayCan standard | CoT |
|---|---|---|---|---|---|---|---|---|---|---|---|
| UL2 | 20B | 34.2 | **51.4** | 59.0 | 53.3 | 13.5 | **14.0** | 57.9 | **65.3** | 20.0 | **41.7** |
| LaMDA | 420M | 20.1 | 19.2 | 46.4 | 24.9 | 1.9 | 1.6 | 50.0 | 49.7 | 7.5 | 7.5 |
| | 2B | 20.2 | 19.6 | 52.6 | 45.2 | 8.0 | 6.8 | 49.3 | 57.5 | 8.3 | 8.3 |
| | 8B | 19.0 | 20.3 | 54.1 | 46.8 | 9.5 | 5.4 | 50.0 | 52.1 | 28.3 | 33.3 |
| | 68B | 37.0 | **44.1** | 59.6 | **62.2** | 15.5 | **18.6** | 55.2 | **77.5** | 35.0 | **42.5** |
| | 137B | 53.6 | **57.9** | 62.4 | **65.4** | 21.5 | **26.8** | 59.5 | **85.8** | 43.3 | **46.6** |
| GPT | 350M | 14.7 | 15.2 | 20.6 | 0.9 | 4.3 | 0.9 | 33.8 | 41.6 | 12.5 | 0.8 |
| | 1.3B | 12.0 | 19.2 | 45.8 | 35.7 | 4.0 | 1.4 | 0.0 | 26.9 | 20.8 | 9.2 |
| | 6.7B | 19.0 | **24.0** | 53.6 | 50.0 | 8.9 | 4.9 | 0.0 | 4.4 | 17.5 | **35.0** |
| | 175B | 79.5 | 73.5 | 65.9 | 65.4 | 43.8 | **52.1** | 69.6 | **82.4** | 81.7 | **87.5** |
| Codex | - | 82.3 | 77.9 | 67.1 | **73.2** | 49.0 | **64.8** | 71.7 | **98.5** | 85.8 | **88.3** |
| PaLM | 8B | 19.8 | **24.9** | 55.6 | 53.5 | 12.9 | 13.1 | 55.1 | **75.2** | 34.2 | **40.0** |
| | 62B | 65.4 | **68.1** | 58.4 | **63.4** | 29.8 | **44.7** | 72.1 | **93.6** | 65.8 | **70.0** |
| | 540B | 78.1 | **79.9** | 68.6 | **77.8** | 49.0 | **65.3** | 80.5 | **95.4** | 80.8 | **91.7** |

## Math Word Problems (free response)

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?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

## Math Word Problems (multiple choice)

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

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. 9 + 90(2) + 401(3) = 1392. The answer is (b).

## CSQA (commonsense)

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

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).

## StrategyQA

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about 0.6 g/cm^3, which is less than water. Thus, a pear would float. So the answer is no.

## Date Understanding

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?

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.

## Sports Understanding

Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."

A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.

## SayCan (Instructing a robot)

Human: How would you bring me something that isn't a fruit?

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.
Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().

## Last Letter Concatenation

Q: Take the last letters of the words in "Lady Gaga" and concatenate them.

A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.

## Coin Flip (state tracking)

Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?

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.

# Results

- In domain and out-of-domain evaluations



Legend:
- Standard prompting
- Chain-of-thought prompting

Letter Concat: 2 (in domain)   Letter Concat: 4 (OOD)

Coin Flip: 2 (in domain)   Coin Flip: 4 (OOD)

Solve rate (%)
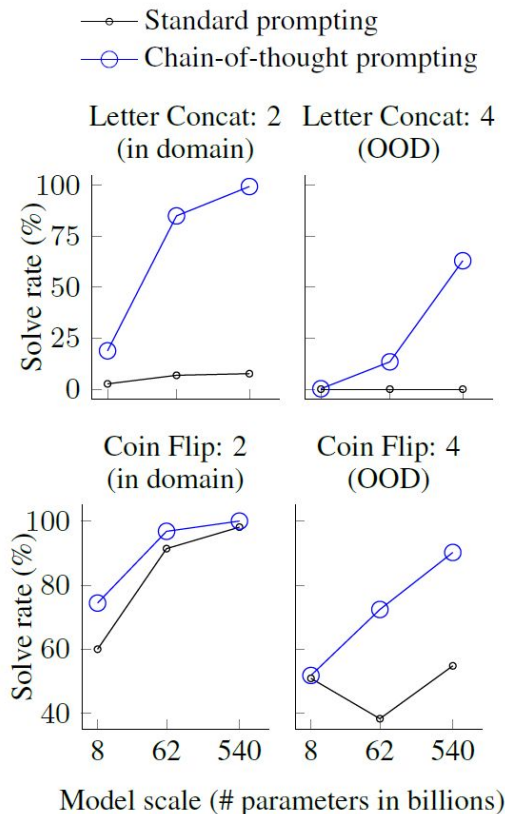
Model scale (# parameters in billions)

Table 5: Standard prompting versus chain of thought prompting enables length generalization to longer inference examples on two symbolic manipulation tasks.

| Model | | Last Letter Concatenation | | | | | | Coin Flip (state tracking) | | | | | |
| | | 2 | | OOD: 3 | | OOD: 4 | | 2 | | OOD: 3 | | OOD: 4 | |
| | | standard | CoT | standard | CoT | standard | CoT | standard | CoT | standard | CoT | standard | CoT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| UL2 | 20B | 0.6 | **18.8** | 0.0 | 0.2 | 0.0 | 0.0 | 70.4 | 67.1 | 51.6 | 52.2 | 48.7 | 50.4 |
| LaMDA | 420M | 0.3 | **1.6** | 0.0 | 0.0 | 0.0 | 0.0 | 52.9 | 49.6 | 50.0 | 50.5 | 49.5 | 49.1 |
| | 2B | 2.3 | **6.0** | 0.0 | 0.0 | 0.0 | 0.0 | 54.9 | **55.3** | 47.4 | 48.7 | 49.8 | 50.2 |
| | 8B | 1.5 | **11.5** | 0.0 | 0.0 | 0.0 | 0.0 | 52.9 | **55.5** | 48.2 | 49.6 | 51.2 | 50.6 |
| | 68B | 4.4 | **52.0** | 0.0 | **0.8** | 0.0 | **2.5** | 56.2 | **83.2** | 50.4 | **69.1** | 50.9 | **59.6** |
| | 137B | 5.8 | **77.5** | 0.0 | **34.4** | 0.0 | **13.5** | 49.0 | **99.6** | 50.7 | **91.0** | 49.1 | **74.5** |
| PaLM | 8B | 2.6 | **18.8** | 0.0 | 0.0 | 0.0 | **0.2** | 60.0 | **74.4** | 47.3 | **57.1** | 50.9 | **51.8** |
| | 62B | 6.8 | **85.0** | 0.0 | **59.6** | 0.0 | **13.4** | 91.4 | **96.8** | 43.9 | **91.0** | 38.3 | **72.4** |
| | 540B | 7.6 | **99.4** | 0.2 | **94.8** | 0.0 | **63.0** | 98.1 | **100.0** | 49.3 | **98.6** | 54.8 | **90.2** |

**Harvard** John A. Paulson
**School of Engineering**
and Applied Sciences

# Limitations & Discussion

Limitations:

- Although chain of thought emulates the thought processes of human reasoners, this does not answer whether the neural network is actually "reasoning.
- The emergence of chain-of-thought reasoning only at large model scales makes it costly to serve in real-world applications.

Discussion:

- What are some potential limitations of using chain-of-thought prompting for large language models? Are there certain types of reasoning tasks that may not be well-suited to this approach?
- What other prompting methods might expand the range of tasks that language models can solve?
- The authors argue that chain-of-thought prompting has the potential to improve the interpretability of language models. Do you agree with this claim and why?

# Thanks & Questions