



Augmented Language Models : A Survey

Based on

Mialon, Gr  goire, et al. "Augmented language models: a survey." <https://arxiv.org/abs/2302.07842>

MS, University of Illinois at Urbana Champaign (UIUC)
Gargi Balasubramaniam

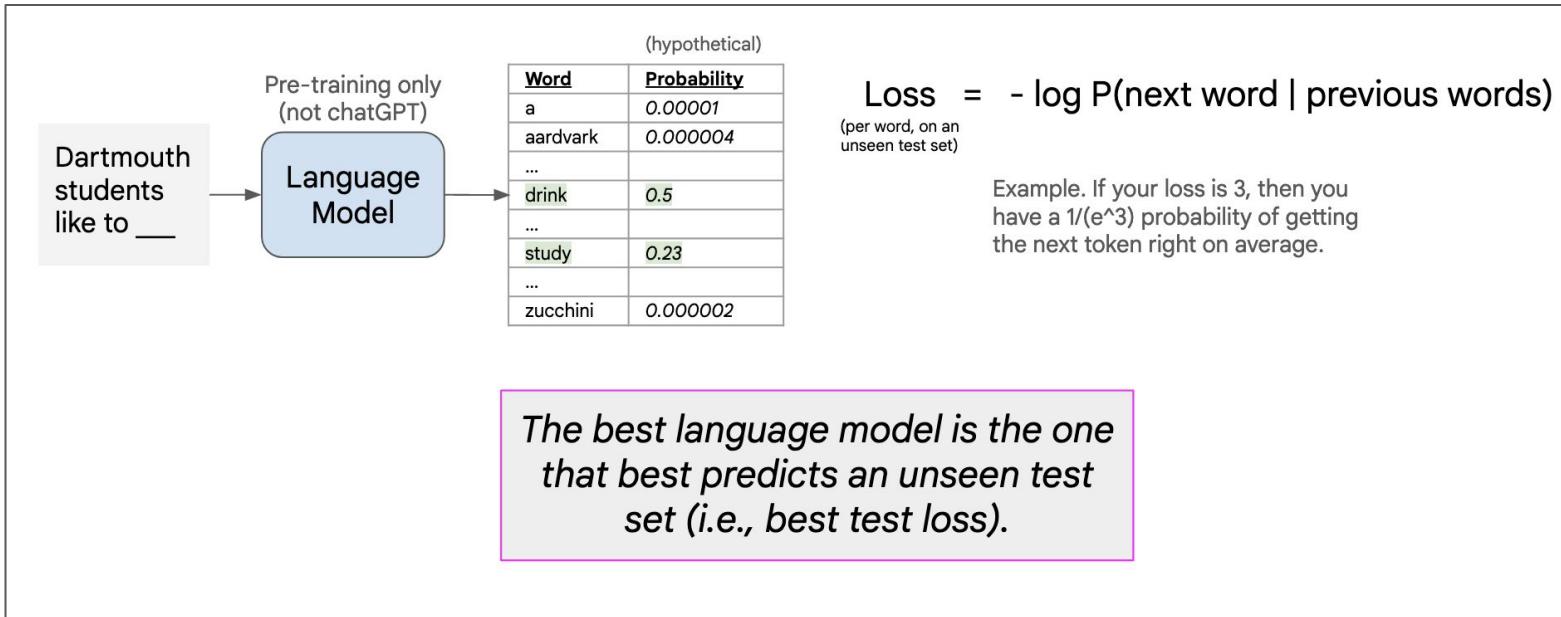
If you find this useful, please consider citing:

Balasubramaniam, Gargi. "Presentation on Augmented Language Models: A Survey." (April 2023).

<https://github.com/bgargi/aug-lang-models-pres>

```
@misc{balasubramaniam23ALM,
author  = {Balasubramaniam, Gargi},
title   = {Presentation on {Augmented Language Models: A Survey}},
year    = {2023},
month   = {April},
howpublished = {PDF slides},
note    = {Based on the survey paper by Mialon, Gr  goire, et al. {"Augmented Language Models: a Survey"} \url{https://arxiv.org/abs/2302.07842}},
url     = {https://github.com/bgargi/aug-lang-models-pres},
}
```

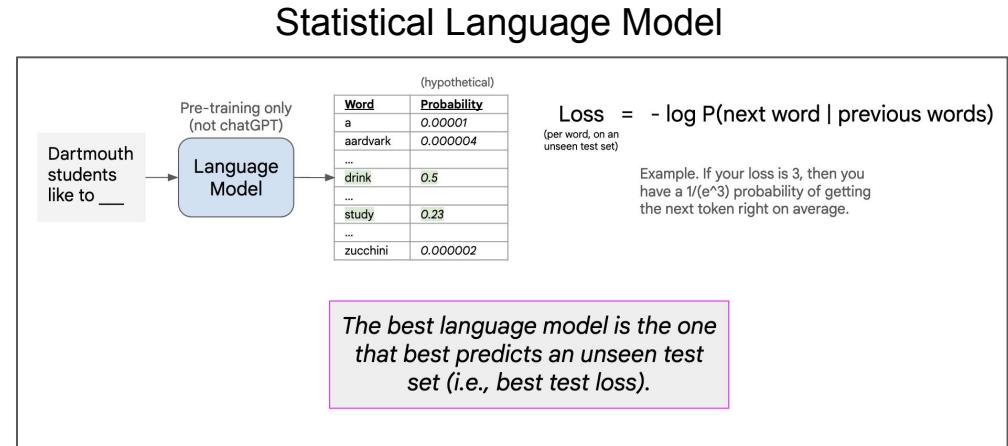
Review



[Source](#)

Why?

Why might one need to augment language models, and how?



Why?

Why might one need to augment language models, and how?

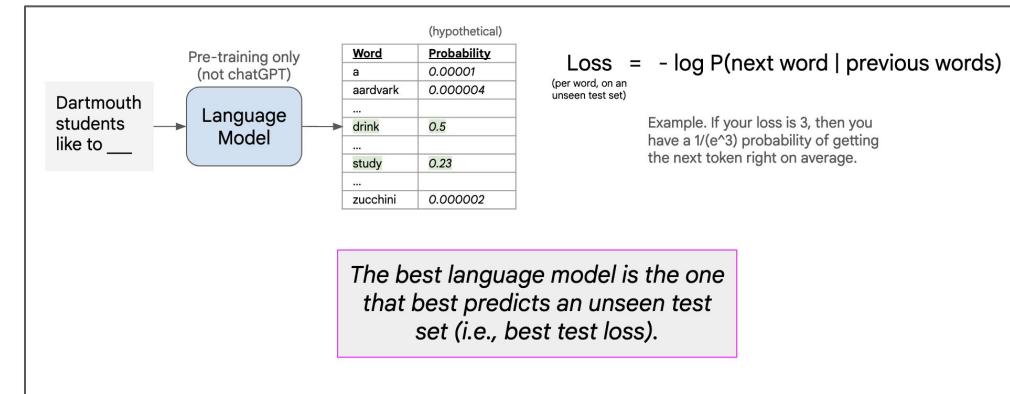
Single Parametric Model



Limited Context



Statistical Language Model



What is possible

What do language models learn from next-word prediction?

| | |
|----------------------------------|--|
| <i>Grammar</i> | In my free time, I like to { <u>run</u> , <u>banana</u> } |
| <i>Lexical semantics</i> | I went to the zoo to see giraffes, lions, and { <u>zebras</u> , <u>spoon</u> } |
| <i>World knowledge</i> | The capital of Denmark is { <u>Copenhagen</u> , <u>London</u> } |
| <i>Sentiment analysis</i> | Movie review: I was engaged and on the edge of my seat the whole time. The movie was { <u>good</u> , <u>bad</u> } |
| <i>Harder sentiment analysis</i> | Movie review: Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was { <u>bad</u> , <u>good</u> } |
| <i>Translation</i> | The word for “pretty” in Spanish is { <u>bonita</u> , <u>hola</u> } |
| <i>Spatial reasoning</i> | [...] Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the { <u>kitchen</u> , <u>store</u> } |
| <i>Math question</i> | First grade arithmetic exam: $3 + 8 + 4 =$ { <u>15</u> , <u>11</u> } |

[thousands (millions?) more]

Extreme multi-task learning!

[Source](#)

.. and what is not

What can't language models learn from next-word prediction?

| | |
|---|--|
| <i>Current world knowledge</i> | The stock price of APPL on March 1st, 2023 is {??? |
| <i>Arbitrarily long arithmetic</i> | $36382894730 + 238302849204 = \text{???}$ |
| <i>Many-step reasoning</i> | Take the nineteenth digit of Pi and multiply it by the e to the fourth power. The resulting ones-digit of the resulting number is {??? |
| <i>Predict the future</i> | The winner of the FIFA world cup in 2026 is {??? |
| <i>Information not in the training data</i> | Jason Wei's favorite color is {??? |
| <i>Extremely long inputs</i> | [2,000 page Harry Potter fan-fiction] What happened after Harry opened the chest for the second time? {??? |

[Source](#)

Why?

Why might one need to augment language models, and **how**?

Statistical Language Model

Single Parametric Model



Limited Context



Today

Induce reasoning

Allow usage of tools

Leverage external data sources (e.g. retrieve documents from a dataset)

Key: (*clever*) context expansion

Disclaimer

- What comes to your mind when you think of “reasoning” in the context of language models?

Disclaimer

- “Reasoning” = decomposing into simpler subtasks and solving those legitimately to obtain a correct answer
 - Simply producing a larger context to correctly produce the missing tokens? Perhaps.

| model | date | developer | # parameters |
|---------------------|----------|----------------------------|-----------------|
| GPT-3 | May 2020 | OpenAI | 175,000,000,000 |
| GPT-Neo | Mar 2021 | EleutherAI | 2,700,000,000 |
| GPT-J | Jun 2021 | EleutherAI | 6,000,000,000 |
| Megatron-Turing NLG | Oct 2021 | Microsoft & Nvidia | 530,000,000,000 |
| Gopher | Dec 2021 | DeepMind | 280,000,000,000 |
| LaMDA | Jan 2022 | Google | 137,000,000,000 |
| GPT-NeoX | Feb 2022 | EleutherAI | 20,000,000,000 |
| Chinchilla | Mar 2022 | DeepMind | 70,000,000,000 |
| PaLM | Apr 2022 | Google | 540,000,000,000 |
| Luminous | Apr 2022 | Aleph Alpha | 70,000,000,000 |
| OPT | May 2022 | Meta | 175,000,000,000 |
| BLOOM | Jul 2022 | Hugging Face collaboration | 175,000,000,000 |
| GPT-3.5 | Nov 2022 | OpenAI | Unknown |
| LLaMA | Feb 2023 | Meta | 65,000,000,000 |
| GPT-4 | Mar 2023 | OpenAI | Unknown |

Not all language models are large.. :)

What does scaling mean and how is it measured?

| | $\alpha \times$ | Model size (# parameters) | \times | Training tokens | = | Training compute |
|--|-----------------|------------------------------|----------|-----------------|---|------------------|
|  BERT-base (2018) | | 109M | | 250B | | 1.6e20 |
|  GPT-3 (2020) | | 175B | | 300B | | 3.1e23 |
|  PaLM (2022) | | 540B | | 780B | | 2.5e24 |

[Source](#)

[Source](#)

Emergence

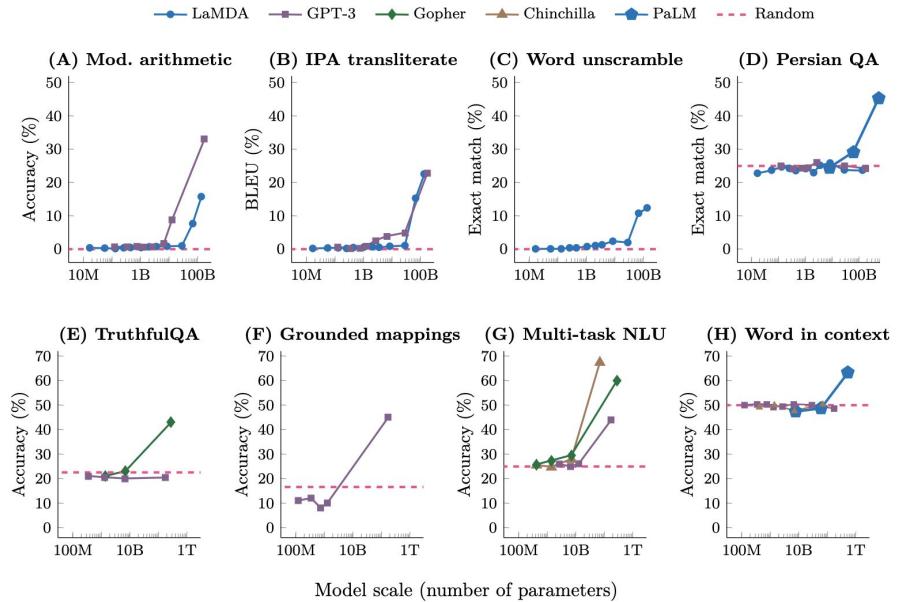
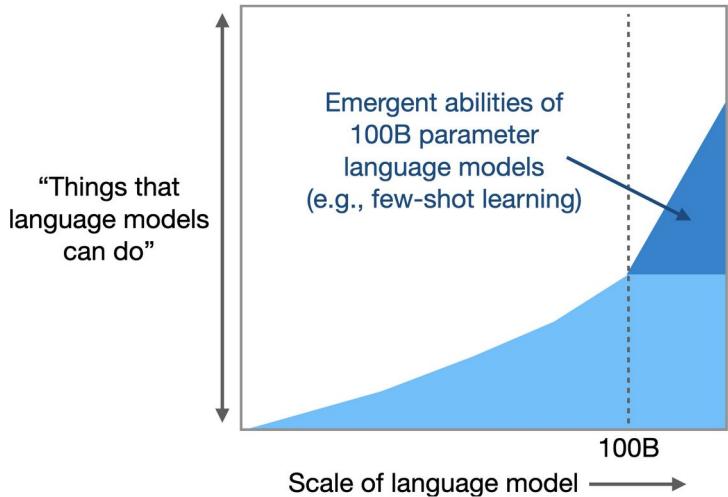


Figure 11: Eight examples of emergence in the few-shot prompting setting. Each point is a separate model.

[Source](#)

[Source](#)

However..

Q1. When was Justin Bieber born?"
(1994)

Q2. "Who was the champion of the Master's tournament in 1994?"

Compositional 2 hop question:

Who was the champion of the Master's Tournament in the year that Justin Bieber was born?

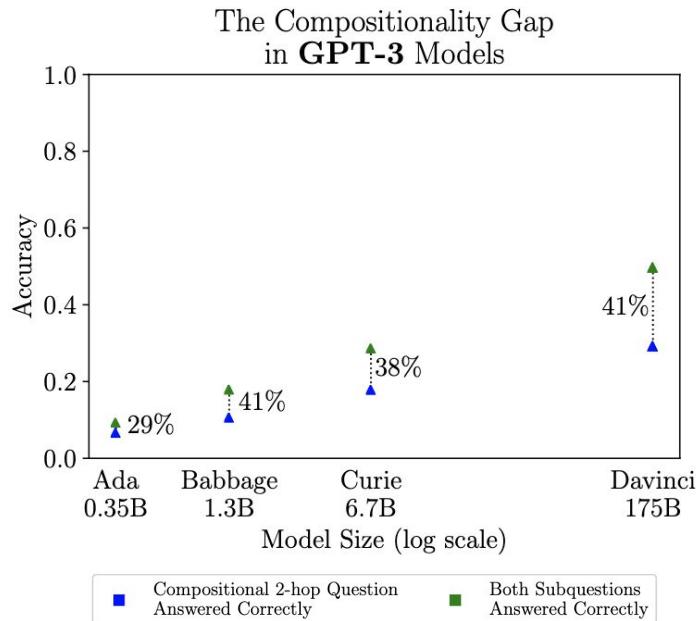
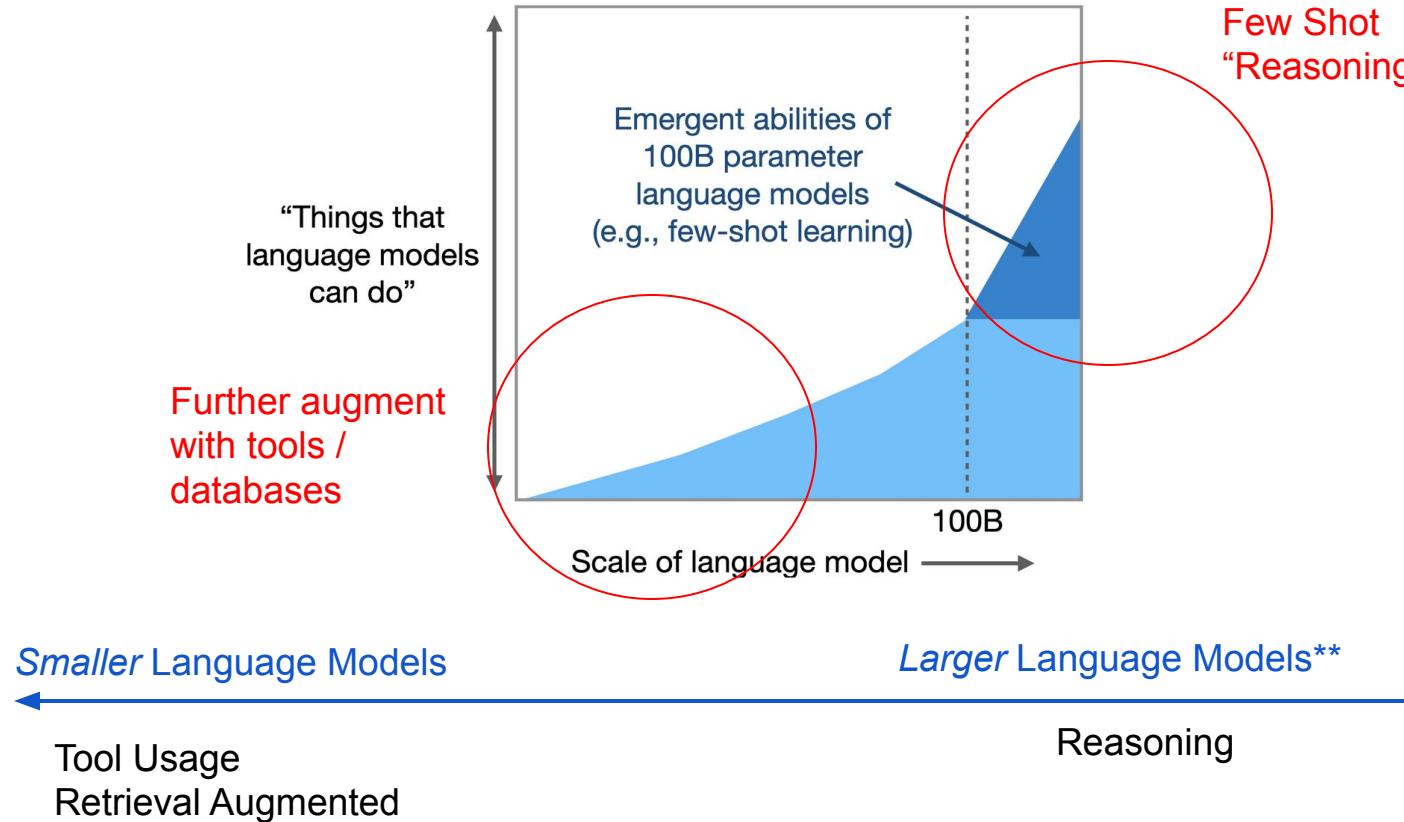


Figure 1: The compositionality gap does not decrease with scale. This plot shows accuracy for the compositional questions (blue) and the two sub-questions that comprise them (green) in the Compositional Celebrities dataset. \triangle are the regular GPT models and \circ and \times are the 001 and 002

Spectrum of ‘augmentation’

My Take



** some exceptions like [UL2](#)

Reasoning



Chain of Thought (CoT) Prompting

Few Shot Demonstrations

Standard Prompting

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

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. 

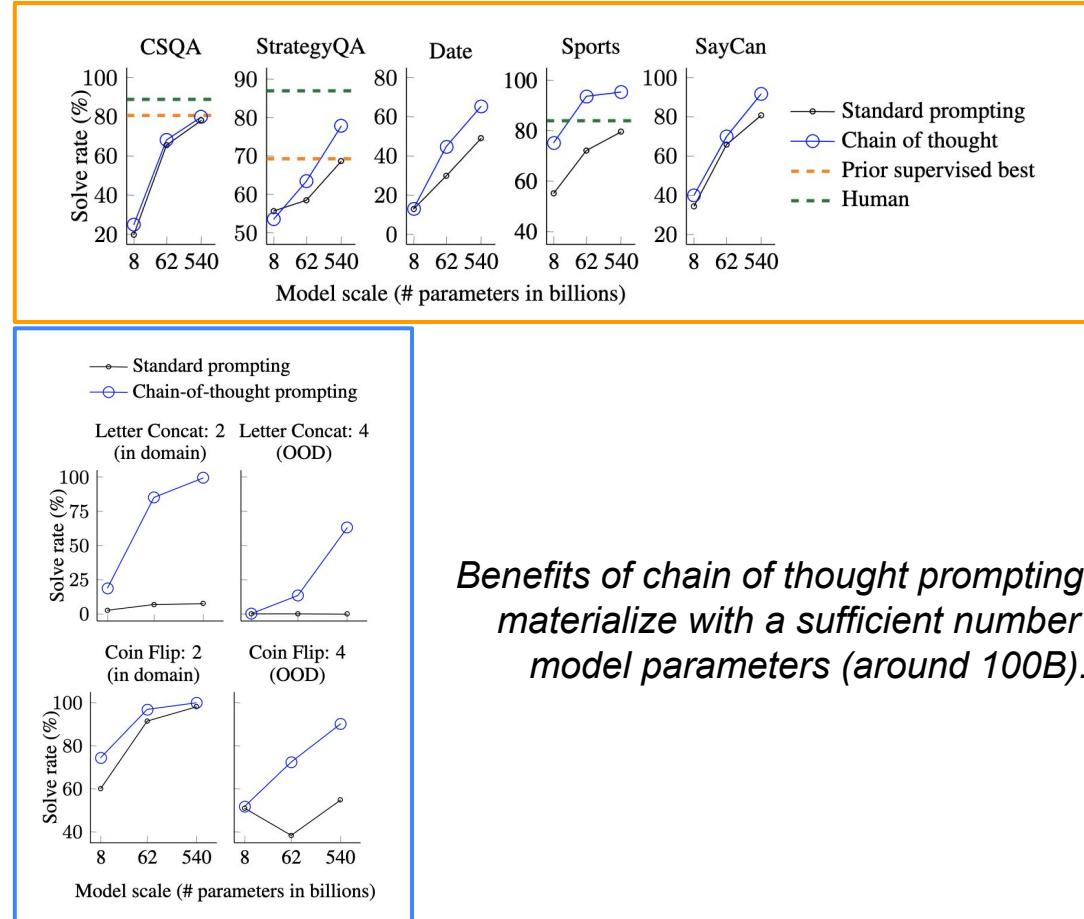
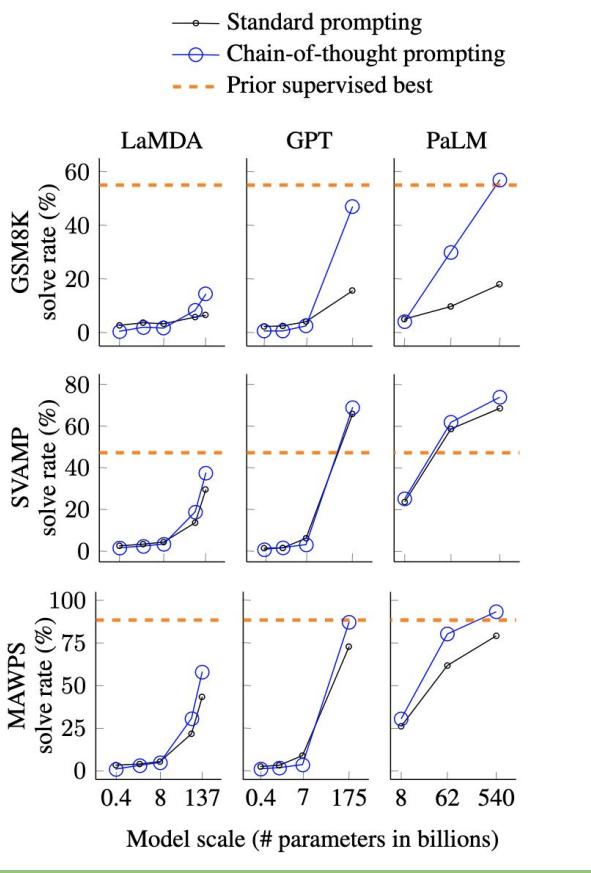
| | | |
|--|--|--|
| <p>Math Word Problems (free response)</p> <p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.</p> | <p>Math Word Problems (multiple choice)</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. $9 + 90(2) + 401(3) = 1392$. The answer is (b).</p> | <p>CSQA (commonsense)</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>StrategyQA</p> <p>Q: Yes or no: Would a pear sink in water?</p> <p>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.</p> | <p>Date Understanding</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>Sports Understanding</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>SayCan (Instructing a robot)</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>Last Letter Concatenation</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>Coin Flip (state tracking)</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> |

+ Arithmetic Reasoning

+ Common Sense Reasoning

+ Symbolic Reasoning

Figure 3: Examples of $\langle \text{input}, \text{chain of thought}, \text{output} \rangle$ triples for arithmetic, commonsense, and symbolic reasoning benchmarks. Chains of thought are highlighted. Full prompts in Appendix G.



Benefits of chain of thought prompting only materialize with a sufficient number of model parameters (around 100B).

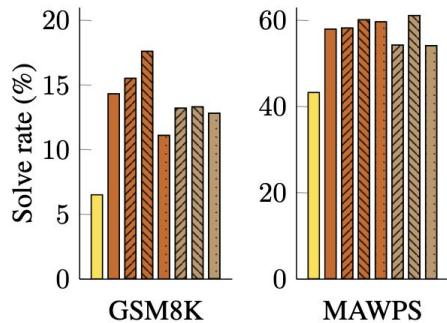
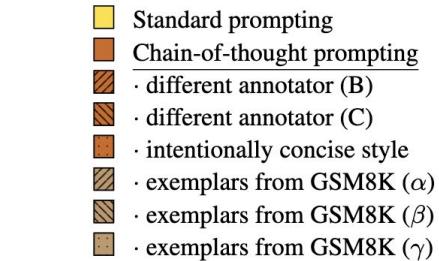


Figure 6: Chain-of-thought prompting has variance for different prompt examples (as expected) but outperforms standard prompting for various annotators as well as for different exemplars.

robustness to

- Annotators
- independently-written chains of thought
- different exemplars

Robustness

Error Analysis on Smaller Models

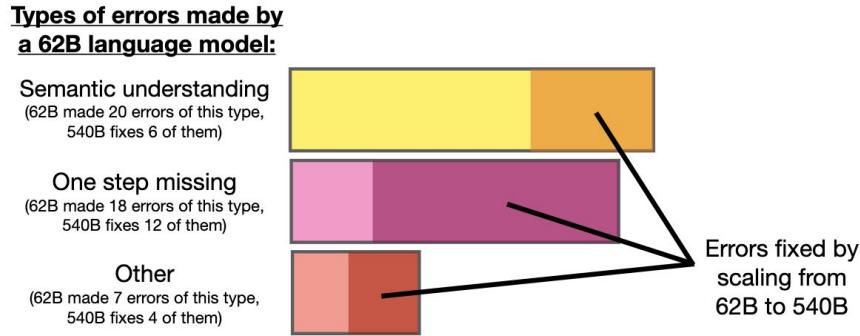


Figure 9: Error analysis of 45 problems that PaLM 62B got incorrect. These errors were categorized into semantic understanding, one step missing, and other. The other category includes hallucinations, repetitive outputs, and symbol mapping errors. Scaling PaLM to 540B fixed a substantial portion of errors in all categories.

Small models:

1. *Even simple operations need larger scale*
2. *Did not generate an answer that could even be parsed*

Chain of Thought (CoT) Prompting

**Zero Shot
Demonstrations**

Chain of Thought (CoT) Prompting

Zero Shot Demonstrations

(a) Few-shot

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(b) Few-shot-CoT

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

(Output) 8 X

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

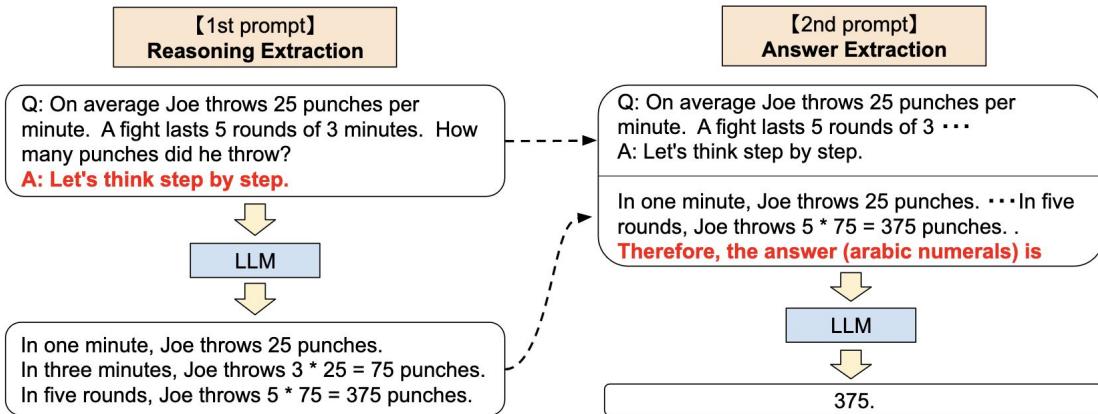
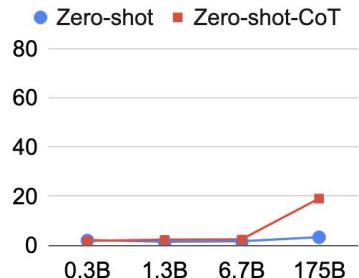
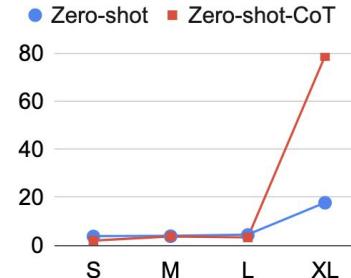


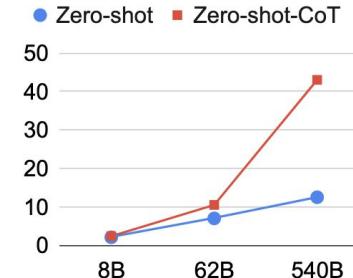
Figure 2: Full pipeline of Zero-shot-CoT as described in § 3: we first use the first “reasoning” prompt to extract a full reasoning path from a language model, and then use the second “answer” prompt to extract the answer in the correct format from the reasoning text.



(a) MultiArith on Original GPT-3



(b) MultiArith on Instruct GPT-3



(c) GMS8K on PaLM

Table 23: Categorization results of produced chain of thought for MultiArith datasets. (*1)
 These categories are cited from Wei et al. [2022].

| Prediction | CoT Category | Zero-Shot-CoT (%) | Few-Shot-CoT (%) |
|------------|-------------------------------|-------------------|------------------|
| Correct | CoT is correct | 94.0 | 98.0 |
| | CoT is incorrect | 6.0 | 2.0 |
| Incorrect | CommonSense Mistake | 10.0 | 23.8 |
| | Factual Mistake | 2.0 | 0.0 |
| | Logical Mistake | 68.0 | 73.8 |
| | - Calculator error (*1) | (8.) | (26.2) |
| | - Symbol mapping error (*1) | (4.) | (2.4) |
| | - One step missing error (*1) | (6.) | (7.1) |
| | - One unnecessary step error | (10.) | (2.4) |
| | - More complicated | (40.) | (35.7) |
| | Others | 20.0 | 2.4 |

Interesting difference in nature of errors

Table 4: Robustness study against template measured on the MultiArith dataset with text-davinci-002.
 (*1) This template is used in Ahn et al. [2022] where a language model is prompted to generate step-by-step actions given a high-level instruction for controlling robotic actions. (*2) This template is used in Reynolds and McDonell [2021] but is not quantitatively evaluated.

| No. | Category | Template | Accuracy |
|-----|-------------|---|-------------|
| 1 | instructive | Let's think step by step. | 78.7 |
| 2 | | First, (*1) | 77.3 |
| 3 | | Let's think about this logically. | 74.5 |
| 4 | | Let's solve this problem by splitting it into steps. (*2) | 72.2 |
| 5 | | Let's be realistic and think step by step. | 70.8 |
| 6 | | Let's think like a detective step by step. | 70.3 |
| 7 | | Let's think | 57.5 |
| 8 | | Before we dive into the answer, | 55.7 |
| 9 | | The answer is after the proof. | 45.7 |
| 10 | misleading | Don't think. Just feel. | 18.8 |
| 11 | | Let's think step by step but reach an incorrect answer. | 18.7 |
| 12 | | Let's count the number of "a" in the question. | 16.7 |
| 13 | | By using the fact that the earth is round, | 9.3 |
| 14 | irrelevant | By the way, I found a good restaurant nearby. | 17.5 |
| 15 | | Abrakadabra! | 15.5 |
| 16 | | It's a beautiful day. | 13.1 |
| - | | (Zero-shot) | 17.7 |

CoT self consistency

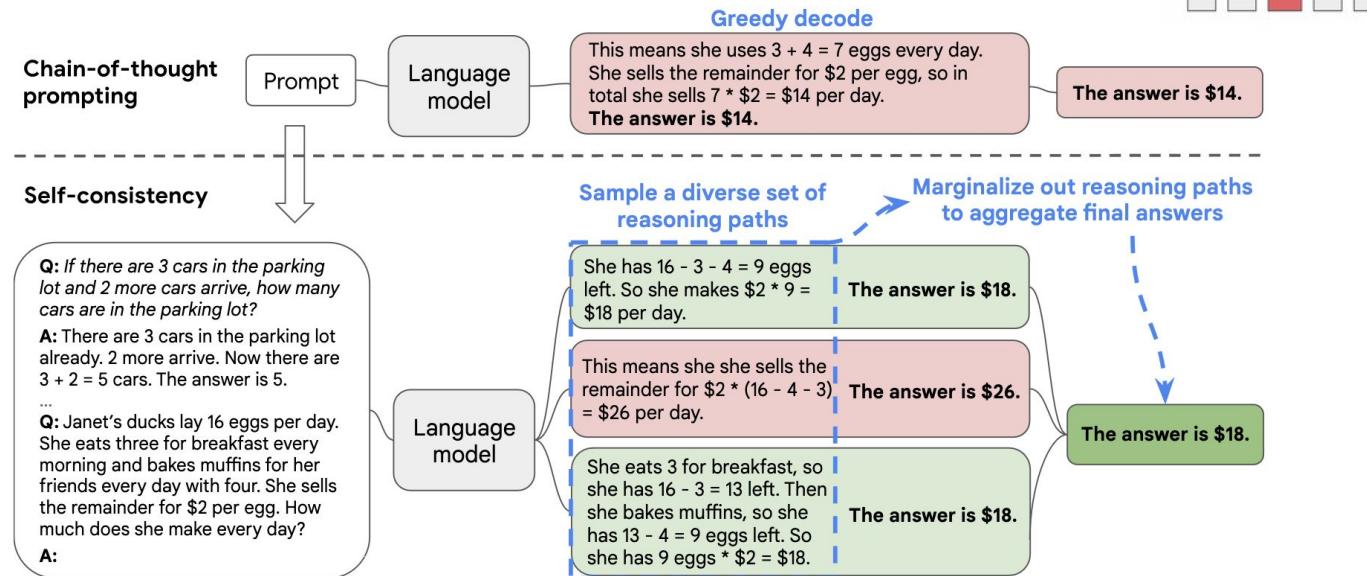
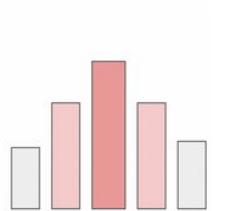


Figure 1: The self-consistency method contains three steps: (1) prompt a language model using chain-of-thought (CoT) prompting; (2) replace the “greedy decode” in CoT prompting by sampling from the language model’s decoder to generate a diverse set of reasoning paths; and (3) marginalize out the reasoning paths and aggregate by choosing the most consistent answer in the final answer set.

[Source](#)



[Image Source](#)

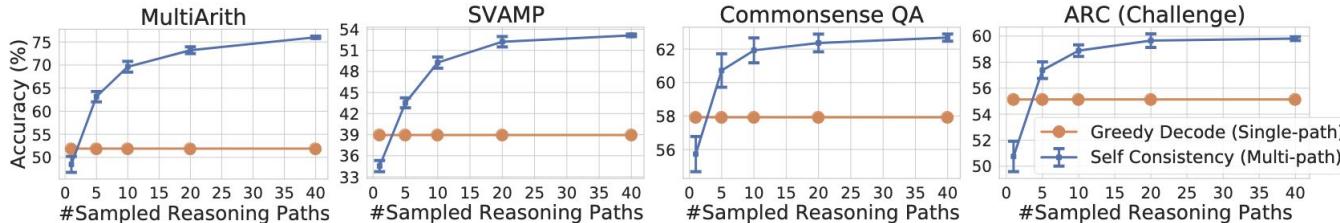


Figure 2: Self-consistency (blue) significantly improves accuracy over CoT-prompts with greedy decoding (orange) across arithmetic and commonsense reasoning tasks, over LaMDA-137B. Sampling a higher number of diverse reasoning paths consistently improves reasoning accuracy.

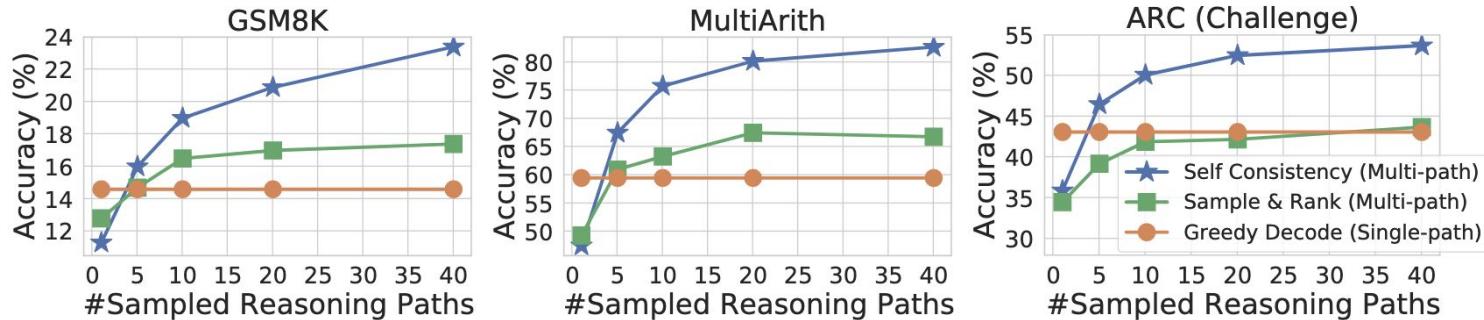


Figure 3: Self-consistency significantly outperforms sample-and-rank with the same # of samples.

Observe the trend of increase. Any comments?

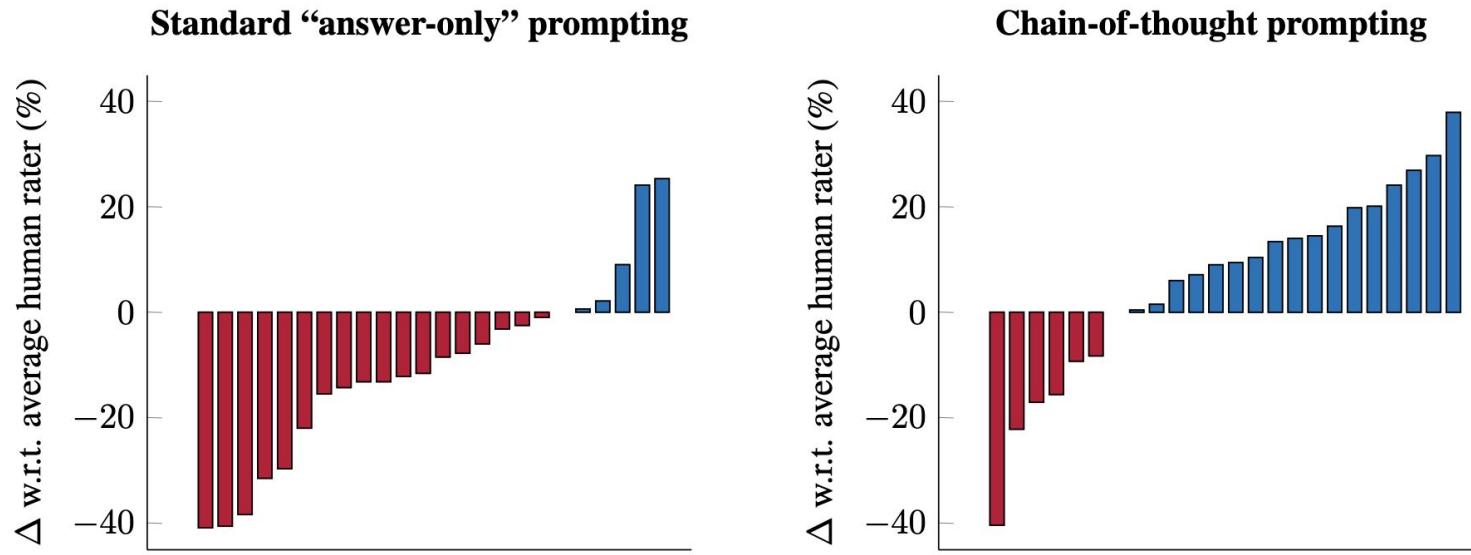
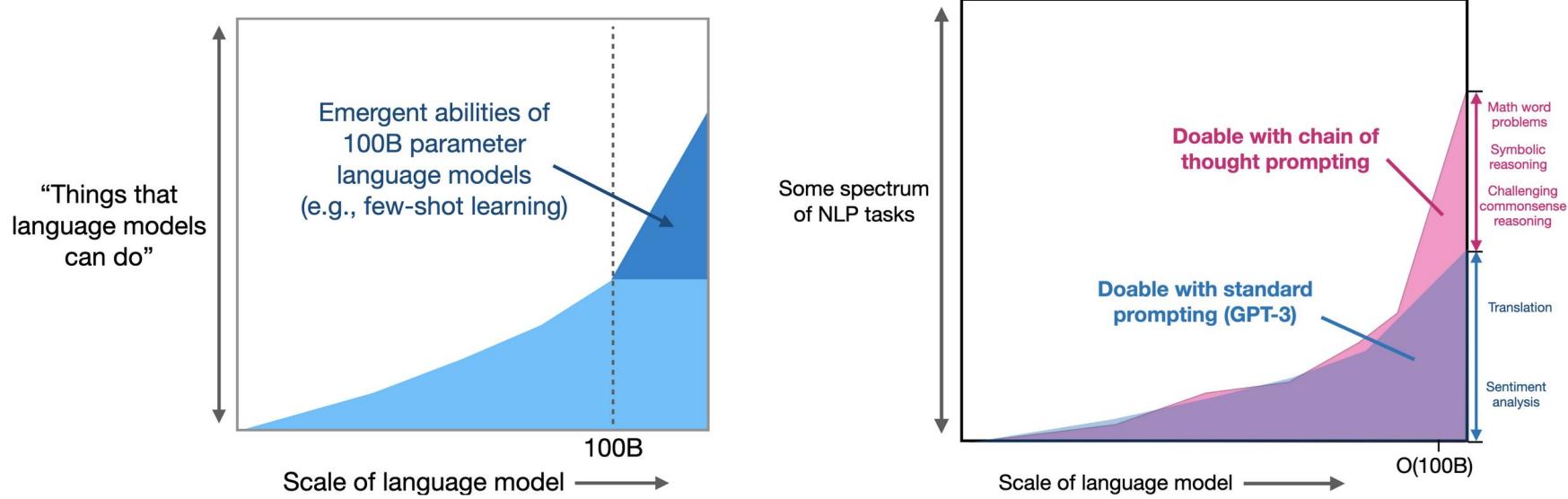


Figure 1: Per-task delta between Codex (code-davinci-002) and the average human-rater performance on 23 challenging tasks in BIG-Bench Hard, for standard “*answer-only*” (left) and *chain-of-thought* (right) prompting.

Closing the compositionality gap?



[Source](#)

However,

CoT may not always be effective for ChatGPT , depends on the task



#ChatGPT spontaneously follows an instruction (#CoT) even when no instruction is given. But its performance drops when explicitly provided with #CoT. This could be a result of instruction memorization/overfitting during instruction finetuning (IFT). Paper: [arxiv.org/pdf/2304.03262...](https://arxiv.org/pdf/2304.03262.pdf)

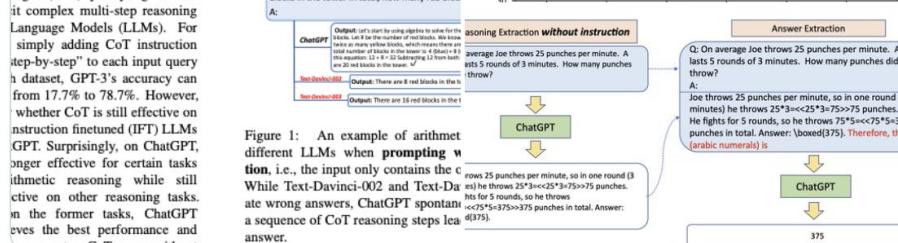
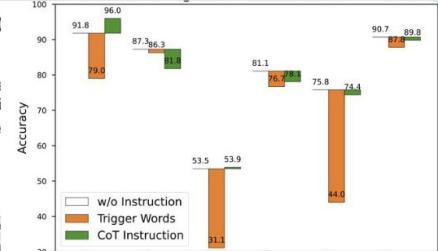
do you need Chain-of-Thought Prompting for ChatGPT

* Lichang Chen* Heng Huang Tianqi Chen
Maryland University of Maryland University of Maryland
id.edu bobchen@umd.edu heng@umd.edu tianqi@umd.edu

Abstract

Chain-of-thought (CoT) prompting can effectively improve complex multi-step reasoning Language Models (LLMs). For example, by simply adding CoT instruction "step-by-step" to each input query dataset, GPT-3's accuracy can improve from 17.7% to 78.7%. However, whether CoT is still effective on instruction finetuned (IFT) LLMs like ChatGPT is surprising. Surprisingly, on ChatGPT, CoT is more effective for certain tasks than arithmetic reasoning while still effective on other reasoning tasks. In the former tasks, ChatGPT achieves the best performance and generates CoT even without

Figure 1: An example of arithmetic reasoning tasks showing different LLMs' performance when prompted with or without instruction. The figure includes a bar chart comparing accuracy across three conditions: w/o Instruction, Trigger Words, and CoT Instruction. It also shows examples of how ChatGPT handles arithmetic reasoning with and without CoT prompting.



2

5

19

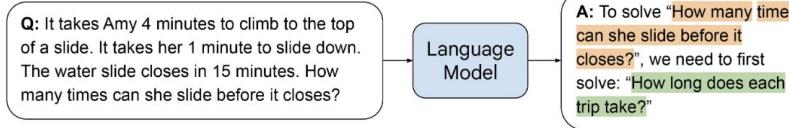
2,131

↑

“Recursive”

Least to most prompting

Stage 1: Decompose Question into Subquestions



Stage 2: Sequentially Solve Subquestions

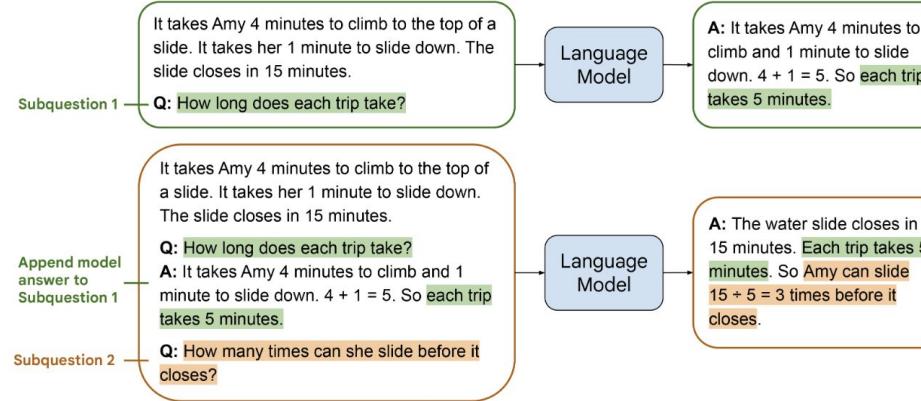
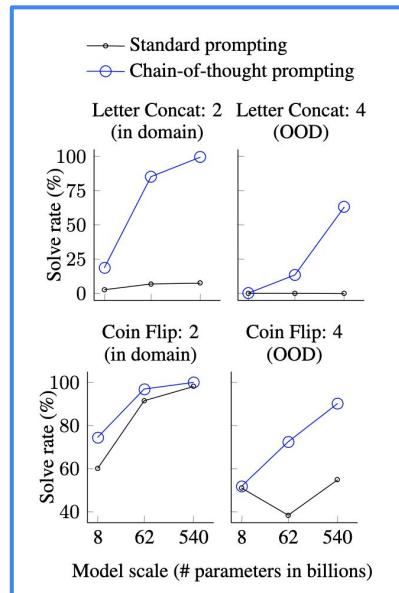


Figure 1: Least-to-most prompting solving a math word problem in two stages: (1) query the language model to decompose the problem into subproblems; (2) query the language model to sequentially solve the subproblems. The answer to the second subproblem is built on the answer to the first subproblem. The demonstration examples for each stage’s prompt are omitted in this illustration.

| | $L = 4$ | $L = 6$ | $L = 8$ | $L = 10$ | $L = 12$ |
|--------------------|-------------|-------------|-------------|-------------|-------------|
| Standard prompting | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| Chain-of-Thought | 84.2 | 69.2 | 50.2 | 39.8 | 31.8 |
| Least-to-Most | 94.0 | 88.4 | 83.0 | 76.4 | 74.0 |

Table 4: Accuracies of different prompting methods on the last-letter-concatenation task. The length of testing lists increases from 4 to 12.



Minerva

| Data source | Proportion of data | Tokens | Present during pretraining |
|-------------------------------|--------------------|--------|----------------------------|
| Math Web Pages | 47.5% | 17.5B | No |
| arXiv | 47.5% | 21.0B | No |
| General Natural Language Data | 5% | >100B | Yes |

PaLM + 118GB dataset of scientific papers using latex, arxiv

Question: A line parallel to $y = 4x + 6$ passes through $(5, 10)$. What is the y -coordinate of the point where this line crosses the y -axis?

Target: Since the line is parallel to $y = 4x + 6$, we know that its slope is 4. Thus the equation of the line is $y = 4x + b$ where b is the y -intercept. Plugging $(5, 10)$ into this gives $10 = 4 \cdot 5 + b \Rightarrow b = [-10]$ which is what we wanted.

Model output: The line is parallel to $y = 4x + 6$, which means that it has the same slope as $y = 4x + 6$. Since the line passes through $(5, 10)$, we can use the point-slope form of a line to find the equation of the line:

$$\begin{aligned}y - 10 &= 4(x - 5) \\y - 10 &= 4x - 20 \\y &= 4x - 10\end{aligned}$$

Now that we have the equation of the line, we can find the y -coordinate of the point where the line crosses the y -axis by substituting $x = 0$ into the equation:

$$y = 4 \cdot 0 - 10 = [-10].$$

Figure 1: Question from the MATH dataset (Algebra 4) with a solution generated by Minerva 62B. The solution was selected by majority vote over 256 samples. This and other MATH questions were prompted with the same 4-shot prompt. The model derives the solution by writing down a line equation, simplifying it, substituting x , and solving for y . The model outputs valid L^AT_EX along with natural language explanations. The approach taken by the model is different from the one used in the ground truth target.

Galactica

- step-by-step reasoning is wrapped with a special token `<work>` and `</work>` to mimic an internal working memory.
- At inference time, the model can be asked explicitly to activate this reasoning mode via the `<work>` token.

Question: A needle 35 mm long rests on a water surface at 20 °C. What force over and above the needle's weight is required to lift the needle from contact with the water surface? $\sigma = 0.0728m$.
`<work>`

$$\sigma = 0.0728N/m$$

$$\sigma = F/L$$

$$0.0728 = F/(2 \times 0.035)$$

$$F = 0.0728(2 \times 0.035)$$

```
calculate.py
"""
f = 0.0728*(2*0.035)
with open("output.txt", "w") as file:
    file.write(str(round(f, 5)))
"""

<run: calculate.py>
```

```
<read: output.txt>
```

```
0.0051
```

```
</work>
```

Answer: $F = 0.0051N$

Figure 4: Working memory example from [Taylor et al. \(2022\)](#). This prompt and its output are seen during LM pre-training.

Scratchpad

DIRECT EXECUTION PREDICTION

Consider the following Python function:

```
def remove_Occ(s,ch):
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    for i in range(len(s) - 1,-1,-1):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    return s
```

Fill in the ??? below:

```
assert remove_Occ("PHP", "P") == ???
```

Large Language Model

```
assert remove_Occ("PHP", "P") == "H"
```

SCRATCHPAD TRACING

Consider the following Python function:

```
def remove_Occ(s,ch):
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    for i in range(len(s) - 1,-1,-1):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    return s
```

```
output = remove_Occ("PHP", "P")
```

What is the execution trace?

Large Language Model

```
state: {}
line: def remove_Occ(s,ch):
state: {"remove_Occ": "<callable_object remove_Occ>"}
line: output = remove_Occ("PHP", "P")
state: {"ch": "P", "s": "PHP"}
line:     for i in range(len(s)):
state: {"ch": "P", "s": "PHP", "i": 0}
line:         if (s[i] == ch):
state: {"ch": "P", "s": "PHP", "i": 0}
line:             s = s[0 : i] + s[i + 1:]
state: {"ch": "P", "s": "HP", "i": 0}
line:                 break
state: {"ch": "P", "s": "HP", "i": 0}
line:     for i in range(len(s) - 1,-1,-1):
state: {"ch": "P", "s": "HP", "i": 1}
line:         if (s[i] == ch):
state: {"ch": "P", "s": "HP", "i": 1}
line:             s = s[0 : i] + s[i + 1:]
state: {"ch": "P", "s": "H", "i": 1}
line:                 break
state: {"ch": "P", "s": "H", "i": 1}
line:             return s
state: {"remove_Occ": "<callable_object remove_Occ>", "output": "H"}
```

Figure 1: Overview of our scratchpad approach applied to predicting code execution and comparison to direct execution prediction. Top: Previous work has shown that large pre-trained models achieve poor performance when asked to directly predict the result of executing given computer code (Austin et al., 2021). Bottom: In this work, we show that training models to use a *scratchpad* and predict the program execution trace line-by-line can lead to large improvements in execution prediction performance. N.B. Although the example above only has one loop iteration for each loop, all loops are unrolled across time.

Algorithmic Prompting

Table 1: Comparison of different prompting strategies studied in this work. The number of \star indicates the level to which each strategy exhibits the given characteristic. In this work, we refer to the basic approach of presenting only input-target pairs with no additional explanation as *few-shot*, and we refer to prompts that provide explicit instructions but no running examples of a task as *instruction-only*. We see that algorithmic prompt includes both qualities of natural language explanation and explicit intermediate computations.

| Prompt strategy | Input-target pairs | Natural language rationale | Intermediate computations | Rationale diversity |
|------------------|---------------------|----------------------------|---------------------------|---------------------|
| Few-shot | $\star \star \star$ | - | - | - |
| Chain-of-thought | $\star \star \star$ | $\star \star \star$ | \star | $\star \star \star$ |
| Scratchpad | $\star \star \star$ | - | $\star \star$ | - |
| Instruction-only | - | $\star \star \star$ | - | $\star \star \star$ |
| Algorithmic | $\star \star \star$ | $\star \star$ | $\star \star \star$ | \star |

A.3 Additional results on two-number addition

This section includes additional details and results for Section 3. In Figure 10, we provide an illustration of different prompting strategies for two-number addition with differing levels of detail in the explanation.

```

Input:
128+367
Target:
<scratch>
1 2 8 + 3 6 7 , C: 0
1 2 + 3 6 , 5 C: 1
1 + 3 , 9 5 C: 0
, 4 9 5 C: 0
4 9 5
</scratch>4 9 5.

```

Scratchpad

```

Input:
128+367
Target:
<scratch>
[1,2,8] has 3 digits.
[3,6,7] has 3 digits.
[1,2,8] + [3,6,7] , C=0 , 8+7+0=15 , A->5 , C->1
[1,2] + [3,6] , A=[5] , C=1 , 2+6+1=9 , A->9 , C->0
[1] + [3] , A=[9,5] , C=0 , 1+3+0=4 , A->4 , C->0
[] + [] , A=[4,9,5] , C=0 , END
</scratch>
4 9 5

```

Detailed scratchpad w/
comma-delim

Problem: 128+367=?
 Explanation:
 Let's think step by step.
 $128+367=128+300+67=428+67=495$.
 The final Answer is 495.

Chain-of-thought

Problem: 128+367=
 ...
 Length of FN is 3. FN=[1,2,8]. Length of SN is 3. SN=[3,6,7]. FN[3]=8. SN[3]=7. C[3]=0. Since 8+7+0=15, 15%10=5. Length of A is 1. Thus A=[5]. Since (15-5)/10=1, C[2]=1.
 Length of FN is 2. FN=[1,2]. Length of SN is 2. SN=[3,6]. FN[2]=2. SN[2]=6. C[2]=1. Since 2+6+1=9, 9%10=9. Length of A is 2. Thus A=[9,5]. Since (9-9)/10=0, C[1]=0.
 ...
 There are no more digits and C[0]=0. Thus the process is complete. The final Answer is [4,9,5].

Algorithmic prompt

Figure 10: Examples of the two-number addition prompt using different techniques.

Reasoning



Access Tools



Internet Augmented LMs

Split into paragraphs of 6 sentences
Rank based on TF-IDF based cosine similarity
between paragraphs and question

Given a question,
clean text is extracted out of 20 URLs returned by
Google, resulting in a set of documents.

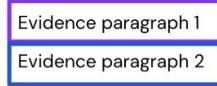
Question:

Who won the NFL football coach of the year?

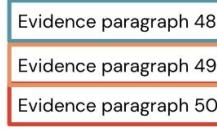
Query Google Search



Chunk documents into passages, rescore them



Most relevant



Least relevant

Internet Augmented LMs

Split into paragraphs of 6 sentences
Rank based on TF-IDF based cosine similarity between paragraphs and question

Given a question,
clean text is extracted out of 20 URLs returned by Google, resulting in a set of documents.

Question: Who won the NFL football coach of the year?

Query Google Search



Chunk documents into passages, rescore them



Evidence paragraph 1
Evidence paragraph 2

Evidence paragraph 48
Evidence paragraph 49
Evidence paragraph 50

Most relevant

Least relevant

Candidate answers for each paragraph

k-shot prompt

Generate answers using LM

Generate scores

Rerank

Combine both: PoE

Evidence: {}
Question: {}
Answer: {}

Evidence: {}
Question: Who won the NFL football coach of the year?
Answer:

Evidence paragraph 1

→ Kevin Stefanski

$p(\text{answer} \mid \text{question}, \text{evidence})$

0.3

$p(\text{question} \mid \text{evidence})$

0.7

PoE

→ 0.5

Evidence paragraph 50

→ Baker Mayfield

0.6

0.1

→ 0.2

"Noisy Channel", maximise:

"RAG"

$$p(a_i \mid q) = \sum_{i=1}^n p_{\text{tf-idf}}(p_i \mid q) \cdot p_{\text{LM}}(a_i \mid q, p_i)$$

$$p(a_i \mid q) = \frac{p_{\text{LM}}(q \mid a_i, p_i) \cdot p_{\text{LM}}(a_i \mid p_i)}{p_{\text{LM}}(q \mid p_i)}$$

Source

| Dataset | SOTA | #sampled answers: 1 | | #sampled answers: 200 | | | |
|------------|-----------|---------------------|--|-----------------------|--------------------|---------------------------------------|-------------------------------------|
| | | CB | OB _{Google} ^{no reranking} | CB | OB _{Gold} | OB _{Google} ^{a q,p} | OB _{Google} ^{PoE} |
| NQ | 51.4 [8] | 21.7 | 23.1 | 25.8 | 61.7 | 32.7 | 38.4 |
| HOTPOTQA | 65.2 [28] | 20.7 | 24.5 | 21.2 | 54.8 | 26.3 | 30.3 |
| FEVER | 73.2 [31] | 44.5 | 52.2 | 44.5 | 66.6 | 52.0 | 57.2 |
| STRATEGYQA | 63.6 [29] | 61.0 | 61.1 | 61.0 | 80.4 | 64.6 | 66.2 |

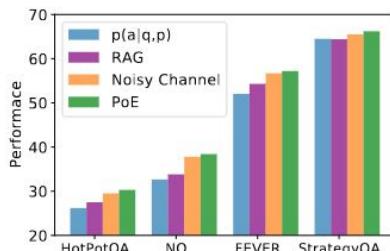


Figure 1: Question answering performance on all 4 datasets using Gopher-280B OB_{Google} model in combination with different answer scoring functions for reranking the answers generated for each of the top 50 Google retrieved paragraphs.

Scaling via inference time, not number of parameters or train time compute

*“searching the Internet is worth
more than 273 billion parameters.”*

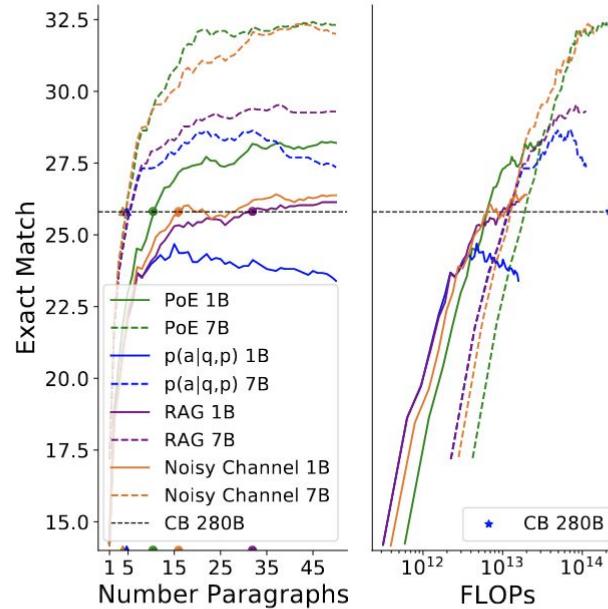


Figure 3: Scaling analysis on the NQ dataset: exact match as a function of number of conditioned paragraphs (left) and FLOPs (right).

Scaling via inference time, not number of parameters or train time compute

“searching the Internet is worth more than 273 billion parameters.”



Natural follow up:
learn to search

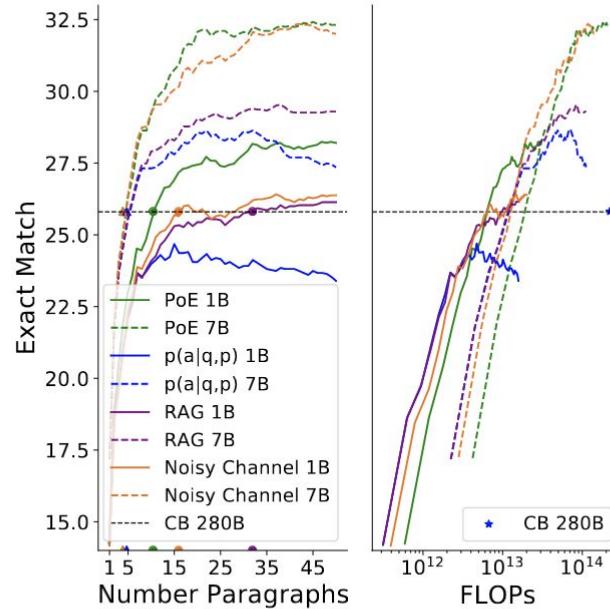


Figure 3: Scaling analysis on the NQ dataset: exact match as a function of number of conditioned paragraphs (left) and FLOPs (right).

SelfAsk

Few shot prompts to demonstrate how / what to ask and then let SelfAsk keep asking / decide when to stop

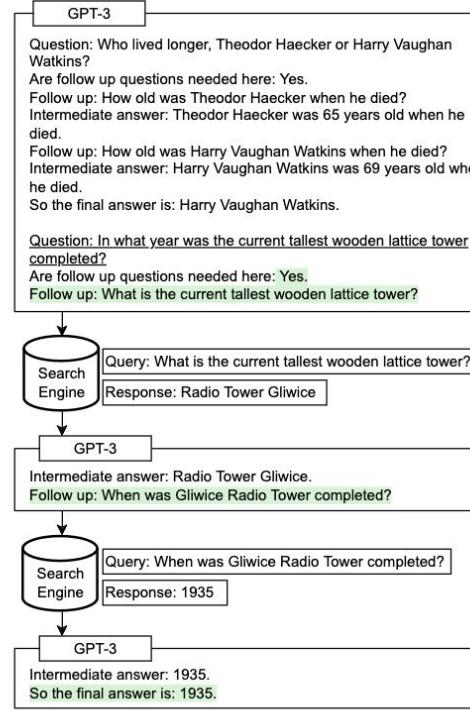


Figure 5: Self-ask + Search Engine: prompt on white background, LM-generated text in green. We start by using a few-shot prompt (reduced here for space) and appending the test-time question (underlined) to it. The LM then generates a follow-up question which we input to an internet search engine. The response is inserted back into the rest of the prompt to let the LM generate the next follow-up question. This process then repeats until the LM decides to output the final answer.

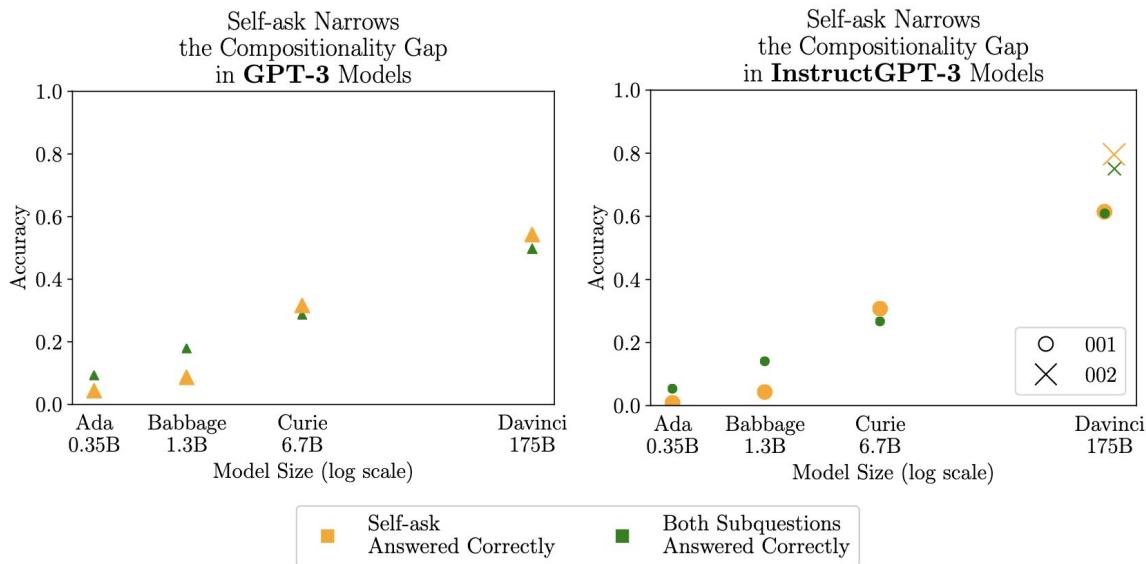


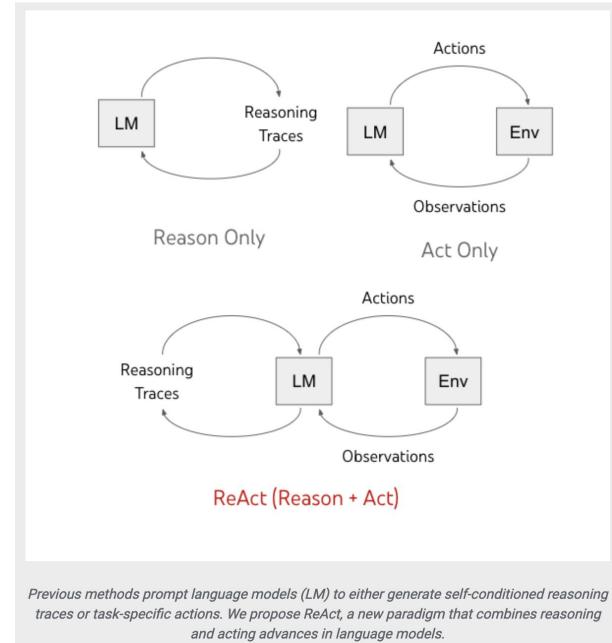
Figure 4: Self-ask is able to narrow and sometimes close the compositionality gap on CC. Here, self-ask uses a 1-shot prompt. Chain of thought performed within 1% of self-ask on this dataset and is not shown here. \triangle are the regular GPT models and \circ and \times are the 001 and 002 instruct models.

ReACT

Synergizing Reasoning and Acting in Language Models

Frozen large language model, PaLM-540B is prompted with few-shot in-context examples to generate

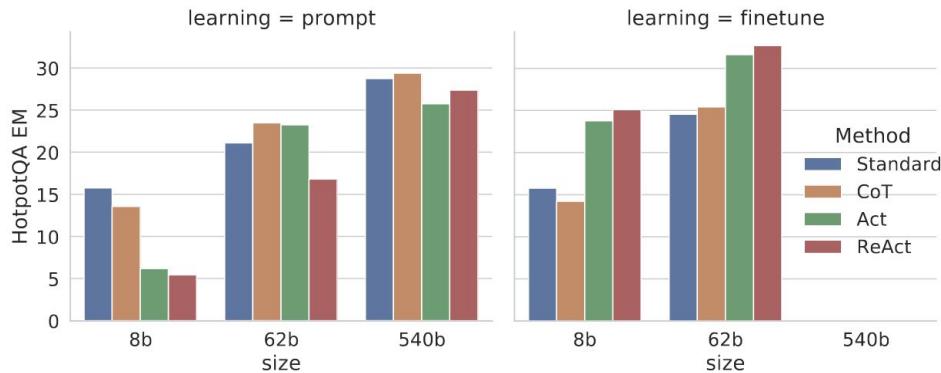
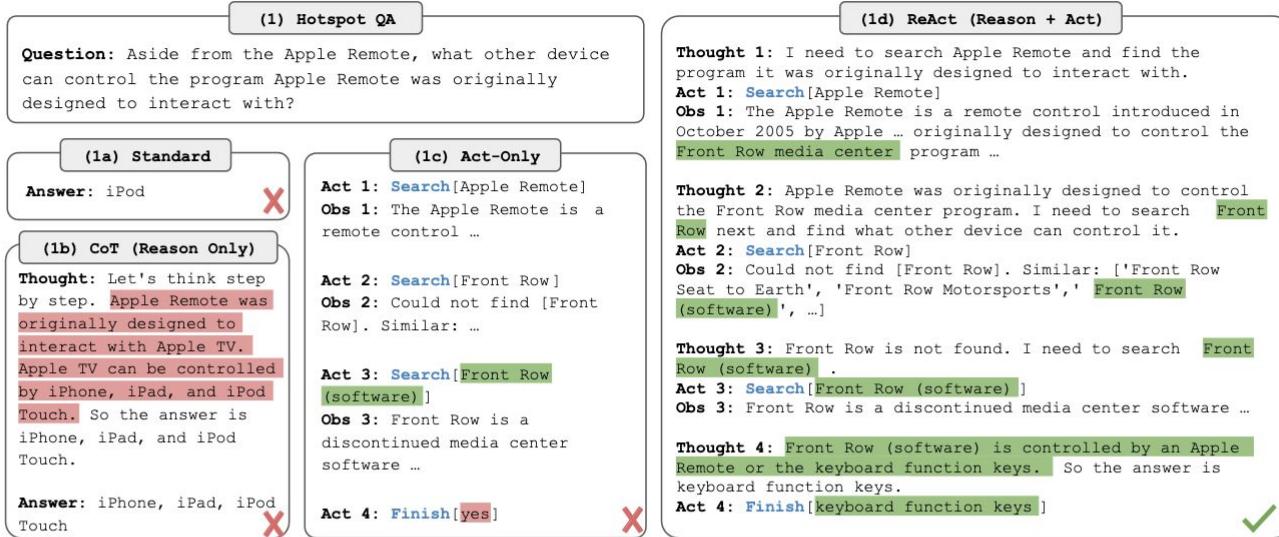
- domain-specific actions and
- free-form language thoughts for task solving
- Each in-context example is a human trajectory of actions, thoughts, and environment observations to solve a task instance



ReAct

| | |
|---------------|---|
| Question | What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into? |
| Thought 1 | I need to search Colorado orogeny, find the area that the eastern sector of the Colorado orogeny extends into, then find the elevation range of the area. |
| Action 1 | Search[Colorado orogeny] |
| Observation 1 | The Colorado orogeny was an episode of mountain building (an orogeny) in Colorado and surrounding areas. |
| Thought 2 | It does not mention the eastern sector. So I need to look up eastern sector. |
| Action 2 | Lookup[eastern sector] |

[Source](#)



PAL

- Main benefit comes from having synergy with interpreter
- <hypothesis> Also trained

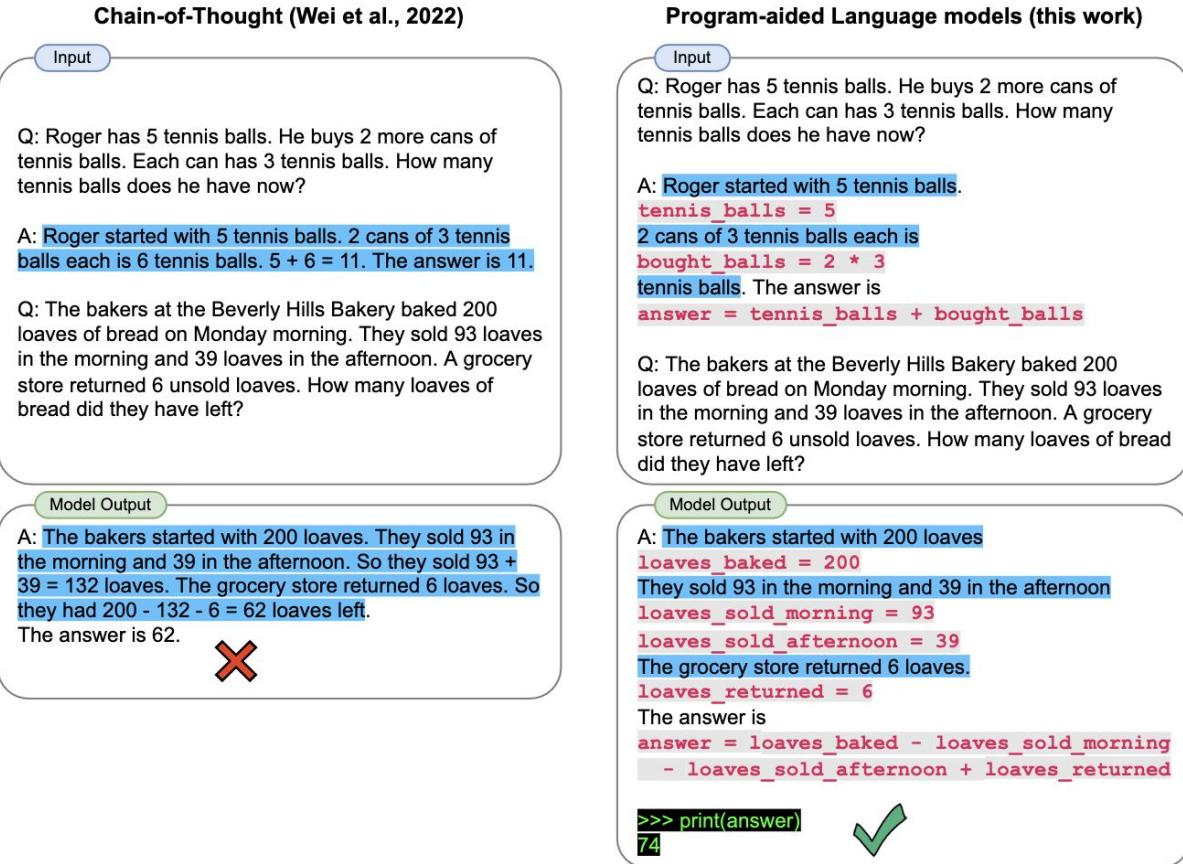


Figure 1: A diagram illustrating PAL: Given a mathematical reasoning question, Chain-of-thought (left) generates intermediate reasoning steps of free-form text. In contrast, Program-aided Language models (PAL, right) generate intermediate steps *and* Python code. This shifts the role of *running* the reasoning steps from the language model to the Python interpreter. The final answer is obtained by running the generated reasoning chain. Chain-of-thought reasoning is highlighted in blue; PAL steps are highlighted in gray and pink; the Python interpreter run is highlighted in black and green.

[Source](#)

Also [POT](#), not covered

every in-context example in PAL is a *pair* $\langle x_i, t_i \rangle$, where $t_j = [s_1, s_2, \dots, s_N]$ with each $s_i \in \text{NL} \cup \text{PL}$, a sequence of tokens in either NL or PL. The complete prompt is thus $p \equiv \langle x_1 \cdot t_1 \rangle \parallel \langle x_2 \cdot t_2 \rangle \parallel \dots \parallel \langle x_k \cdot t_k \rangle$.

Given a test instance x_{test} , we append it to the prompt, and $p \parallel x_{test}$ is fed to the LM. We let the LM generate a prediction t_{test} , which contains both the intermediate steps and their corresponding programmatic statements.

A: Roger started with 5 tennis balls.
`tennis_balls = 5`
2 cans of 3 tennis balls each is
`bought_balls = 2 * 3`
tennis balls. The answer is
`answer = tennis_balls + bought_balls`

Figure 2: A close-up of a single example from a PAL prompt. Chain-of-thought reasoning is highlighted in blue, and PAL programmatic steps are highlighted in gray and pink.

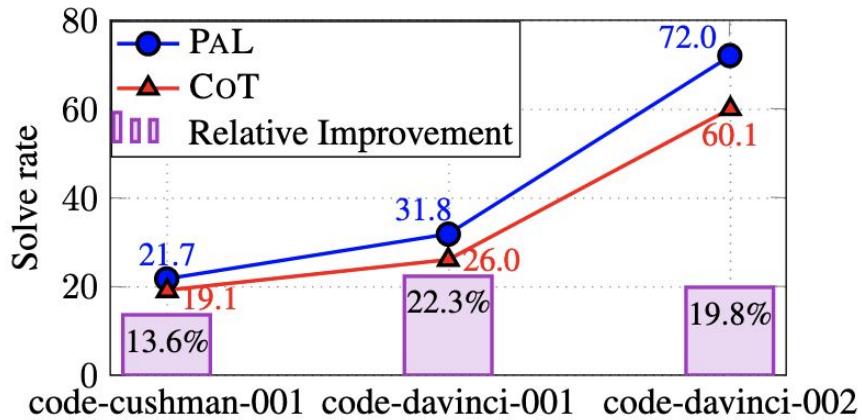


Figure 7: PAL with different models on GSM8K: though the absolute accuracies with code-cushman-001 and code-davinci-001 are lower than code-davinci-002, the relative improvement of PAL over CoT is consistent across models.

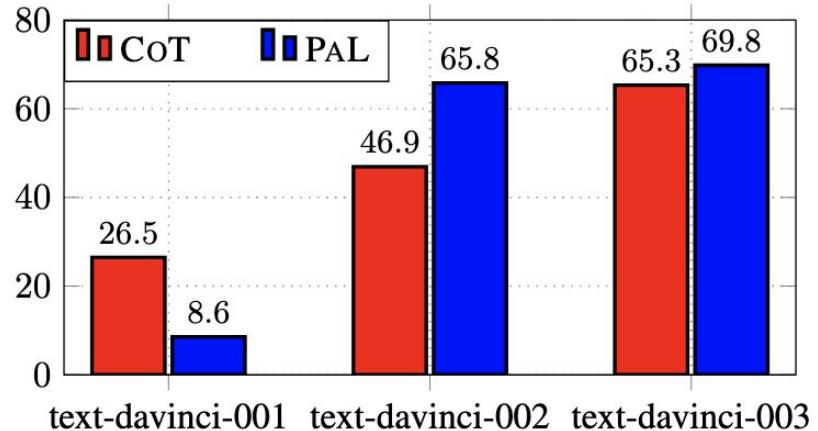


Figure 8: PAL with NL LMs on GSM8K: though CoT outperforms PAL with text-davinci-001, once the base LM is sufficiently strong, PAL is beneficial with text-davinci-002 and text-davinci-003 as well. That is, PAL is not limited to code-LMs only.

Note: No evaluation on smaller models, all are (100B+)

PAL: Program-aided Language Models

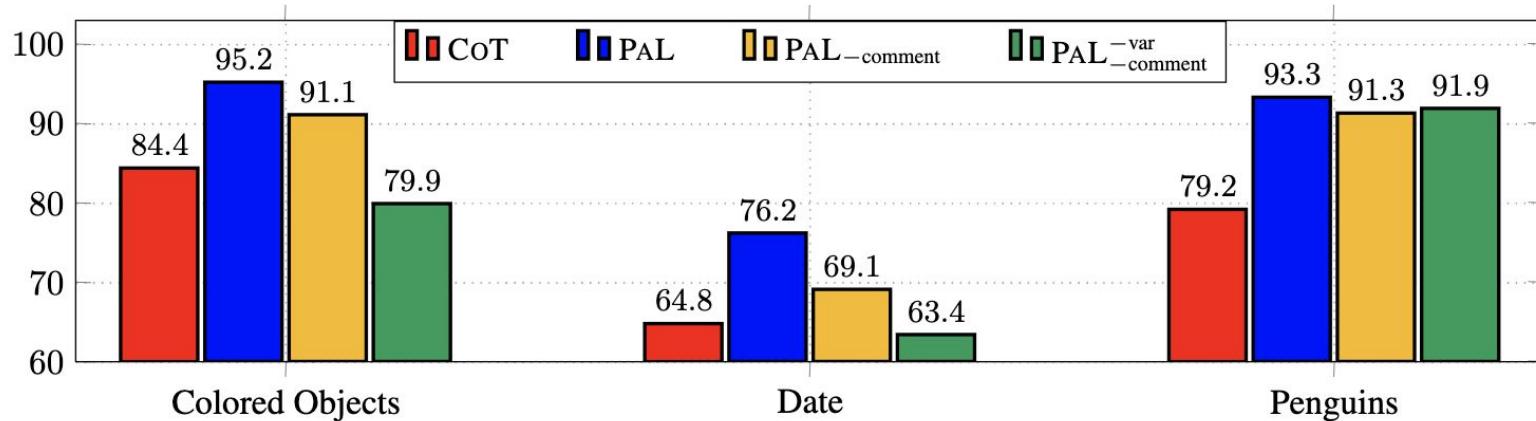


Figure 9: Ablation study of PAL prompt formats. We consider the original PAL prompt, it with natural language comments removed ($\text{PAL}_{-\text{comment}}$), and further variable names replaced with random character ($\text{PAL}_{-\text{var}}_{-\text{comment}}$). As a reference, we also show the CoT performance (blue).

Also works with Least to Most, not only CoT

I. Generalization of PAL to Least-to-Most Prompting

Q: Four years ago, Kody was only half as old as Mohamed. If Mohamed is currently twice 30 years old, how old is Kody?
A: To answer the question "How old is Kody?", we need to know: "How old is Mohamed?", "How old was Mohamed four years ago?", "How old was Kody four years ago?".

(a) Least-to-Most Math Reducing Prompt

Four years ago, Kody was only half as old as Mohamed. If Mohamed is currently twice 30 years old, how old is Kody?

Q: How old was Mohamed four years ago?
A: We were told that Mohamed is currently twice 30 years old, so he is currently $30 * 2 = 60$ years old. That means that four years ago he must have been $60 - 4 = 56$ years old. The answer is 56.

Q: How old is Kody?
A: Four years ago, Kody was half as old as Mohamed, so Kody must have been $56 / 2 = 28$ years old then. Since Kody was 28 years old four years ago, she must now be $28 + 4 = 32$ years old. The answer is 32.

(b) Least-to-Most Math Solving Prompt

```
# Four years ago, Kody was only half as old as Mohamed. If Mohamed is currently twice 30 years old, how old is Kody?
```

```
# How old was Mohamed four years ago?  
mohamed_age_current = 30 * 2  
mohamed_age_4_years_ago = mohamed_age_current - 4
```

```
# Final Question: How old is Kody?  
kody_age_4_years_ago = mohamed_age_4_years_ago / 2  
kody_age_current = kody_age_4_years_ago + 4  
answer = kody_age_current
```

(c) PAL Math Solving Prompt

| Dataset (500 examples) | LEAST-TO-MOST | LEAST-TO-MOST + PAL |
|------------------------|---------------|---------------------|
| GSM8K | 67.2 | 72.8 |
| SVAMP | 75.2 | 78.2 |

Figure 15: Prompts for Math data sets.

Conditioning on
code really helps!

| Model | Accuracy (%) |
|--|--------------|
| OpenAI (<code>text-davinci-002</code>) ^[1] | 15.6 |
| OpenAI (<code>text-davinci-002</code>) + CoT ^[1] | 46.9 |
| OpenAI (<code>text-davinci-002</code>) + CoT + Calculator ^[1] | 46.9 |
| OpenAI (<code>code-davinci-002</code>) ^[1] | 19.7 |
| OpenAI (<code>code-davinci-002</code>) + CoT ^[1] | 63.1 |
| OpenAI (<code>code-davinci-002</code>) + CoT + Calculator ^[1] | 65.4 |
| GPT-3 175B + FT + CoT + Calculator ^[2] | 34.0 |
| GPT-3 175B + FT + CoT + Calculator + Verifier ^[2] | 55.0 |
| PaLM 540B ^[3] | 17.0 |
| PaLM 540B+CoT ^[3] | 54.0 |
| PaLM 540B+CoT+Calculator ^[3] | 58.0 |
| PAL ^[4] | 72.0 |

Table 1: Evaluation of different reasoning methods on GSM8K, a popular reasoning benchmark. FT denotes fine-tuning and CoT denotes chain-of-thought. The reported accuracies are based on [1]: (Wei et al., 2022c); [2]: (Cobbe et al., 2021); [3]: (Chowdhery et al., 2022); and [4]: (Gao et al., 2022).

Other kinds of “code conditioning”

Mindseye

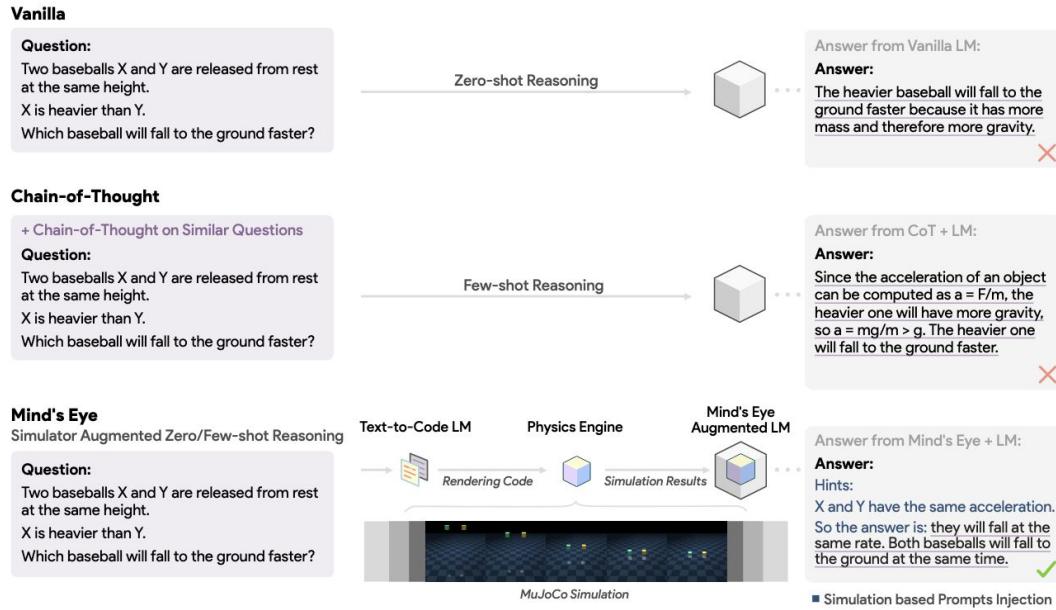
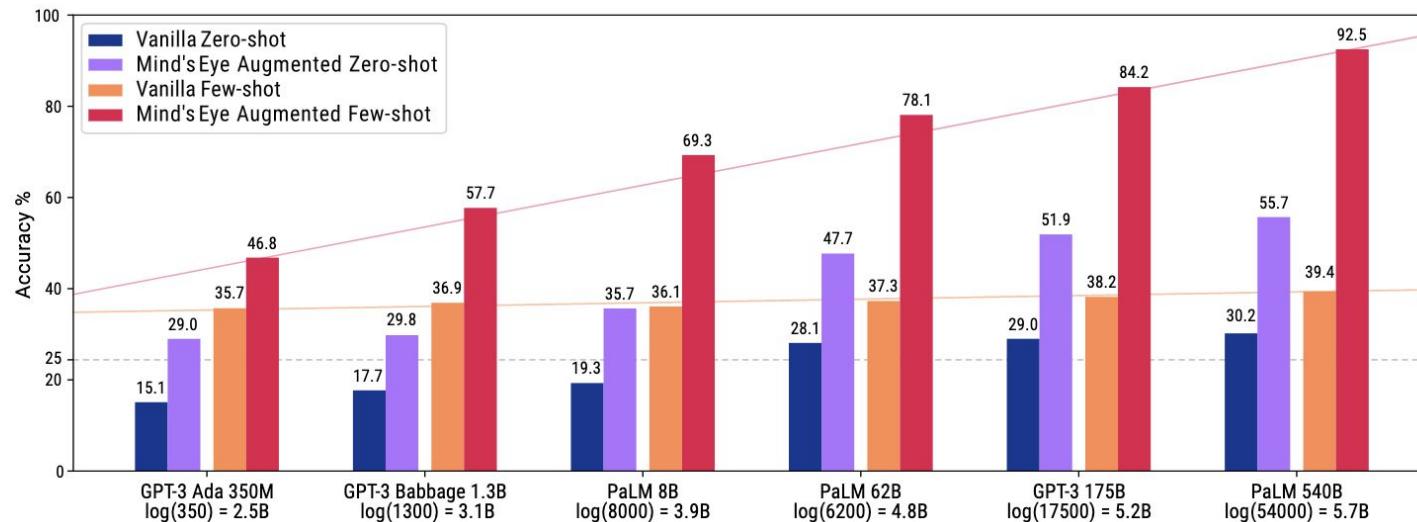


Figure 1: Current language models are still challenged by simple questions that require a good understanding of the physical world. The answer elicited by Chain-of-Thought can still be wrong if the required knowledge is missing or misrepresented in LMs. Mind’s Eye, instead, enables grounded LM reasoning by directly simulating the scene in the given question. Then the LM can reason over the injected ground-truth rationale to generate the correct answers.

Table 2: Comparison in formulation of how the LM inference process is grounded (given question x generating answer y). With the aid of a simulator (i.e., MuJoCo), Mind’s Eye is not only scalable (not requiring human annotation) but also well-grounded with the physical world.

| Method Name | Formulation | Scalable? | Grounded? |
|-------------------------|---|-----------|-----------|
| Zero-shot Reasoner | $y \leftarrow \arg \max_{\hat{y}} \text{LM}(\hat{y} x, \text{"Let's think step by step"})$ | ✓ | ✗ |
| Chain-of-Thought | $y \leftarrow \arg \max_{\hat{y}} \text{LM}(\hat{y} x, \text{Chain-of-Thought} \sim p_{\text{Human}})$ | ✗ | ? |
| Ours: Mind’s Eye | $y \leftarrow \arg \max_{\hat{y}} \text{LM}(\hat{y} x, \text{Simulator}(x, \hat{y}) \sim p_{\text{World}})$ | ✓ | ✓ |



toolformer

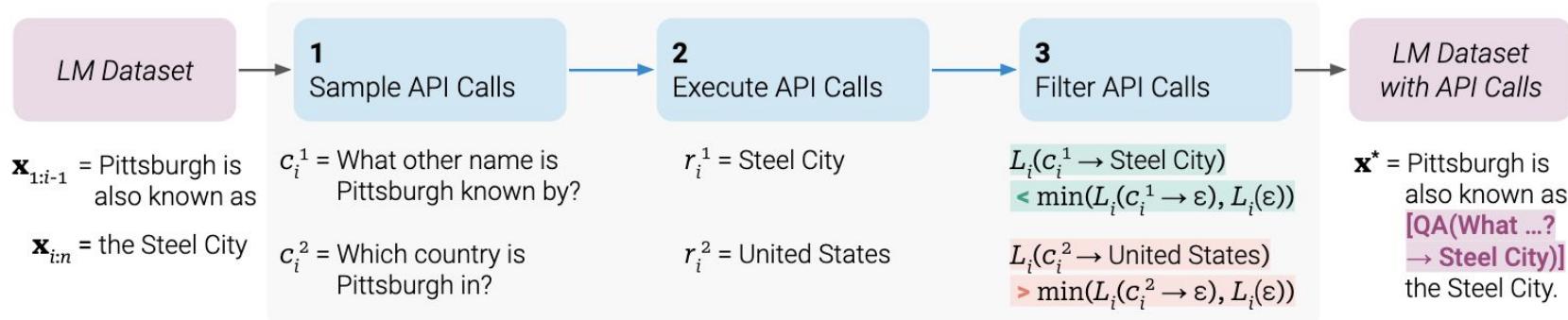


Figure 2: Key steps in our approach, illustrated for a *question answering* tool: Given an input text \mathbf{x} , we first sample a position i and corresponding API call candidates $c_i^1, c_i^2, \dots, c_i^k$. We then execute these API calls and filter out all calls which do not reduce the loss L_i over the next tokens. All remaining API calls are interleaved with the original text, resulting in a new text \mathbf{x}^* .

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

[Source](#)

1. Demonstrate via few shot learning to inject possible API calls
 2. Sample and execute these API calls, keep a top few based on whether they help predict
2. Filter annotations based on whether API calls help model predict future tokens. Use a self-supervised loss to decide which API calls are actually helpful.
- Execute each API call c_i to get corresponding result r_i .
 - Compute weighted cross entropy loss for the LM over tokens x_i, \dots, x_n when the model is prefixed with the prompt. Two versions are computed, one with API result and the other with empty sequence ε .

$$L_i^+ = L_i(e(c_i, r_i))$$
$$L_i^- = \min(L_i(\varepsilon), L_i(e(c_i, \varepsilon)))$$

Only API calls with $L_i^- - L_i^+$ larger than a threshold are kept, meaning that adding this API call and its results help the model predict future tokens.

[Source](#)

Reasoning



Retrieval



Access Tools



IRCoT

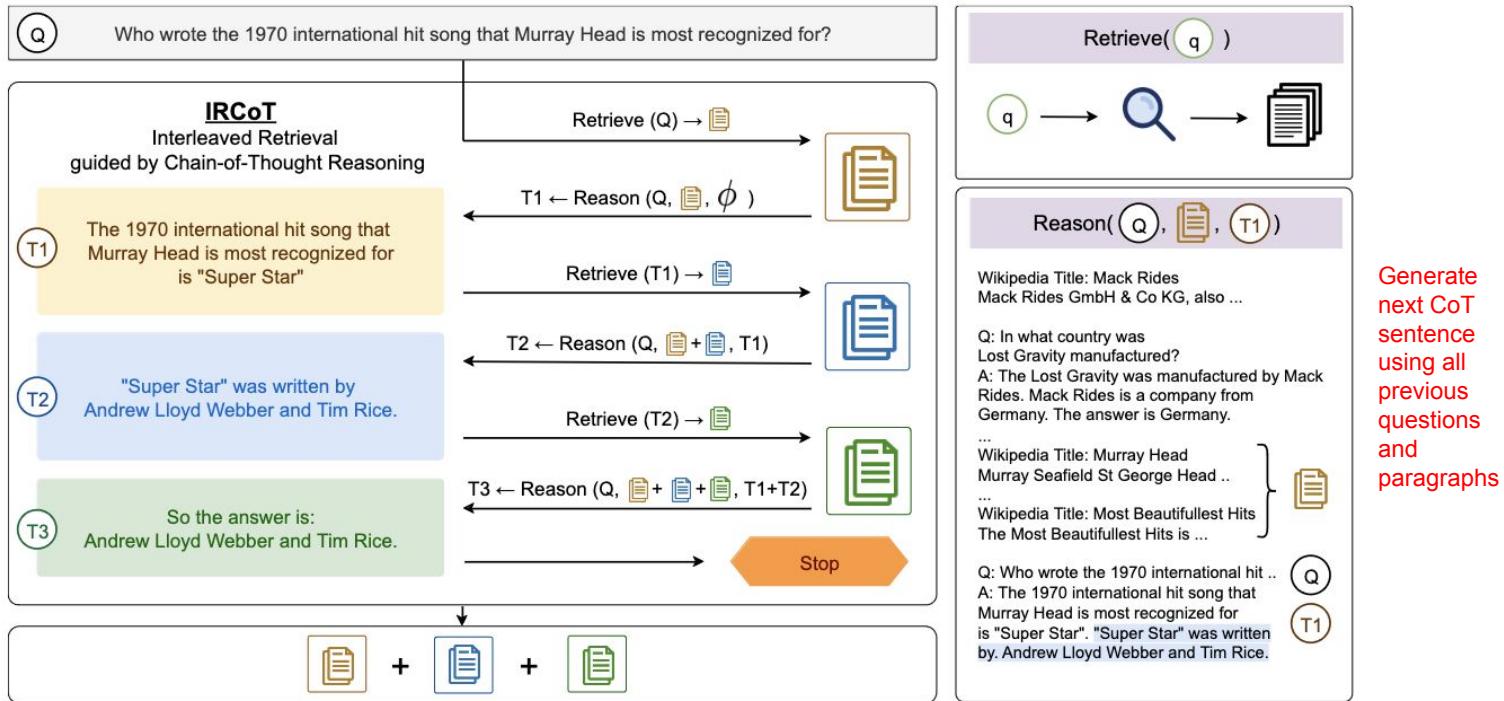


Figure 2: IRCoT interleaves chain-of-thought (CoT) generation and retrieval steps to guide the retrieval by CoT and vice-versa. We start by retrieving K documents using the question as they query and repeat two steps alternately until termination. (i) reason-step generates next CoT sentence based on the question, so far retrieved paragraphs, and CoT sentences. (ii) retrieve-step retrieves K more paragraphs based on the last CoT sentence. The process terminates when the generated CoT has “answer is” or the number of steps exceeds a threshold. The collection of all paragraphs is returned as the retrieval result on the termination.

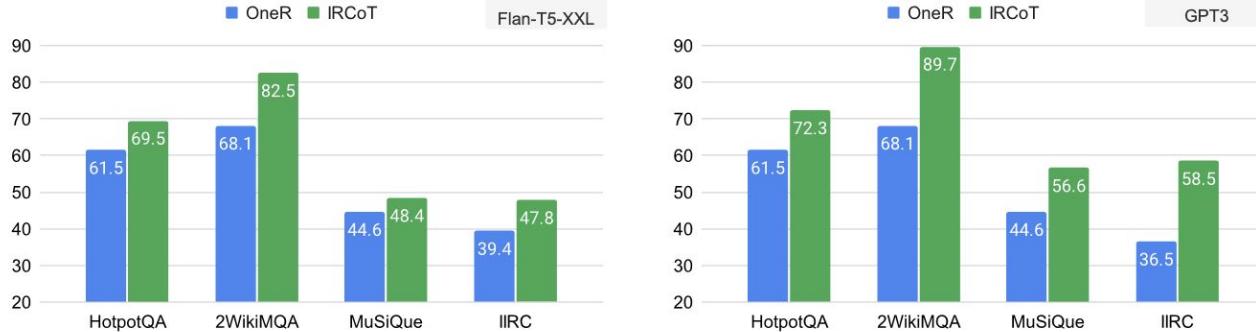


Figure 3: Retrieval recall for one-step retriever (OneR) and IRCoT instantiated from Flan-T5-XXL (left) and GPT3 (right) models. IRCoT outperforms OneR for both models and all datasets. IRCoT with 3B model even outperforms OneR with 58X larger GPT3 model showing the value of improved retrieval.



Figure 5: Retrieval recall for OneR (bottom line) and IRCoT (top line) for language models of increasing sizes: Flan-T5-base (0.3B), Flan-T5-large (0.7B), Flan-T5-XL (3B), Flan-T5-XXL (11B), GPT3 (175B) on HotpotQA (left), 2WikiMultihopQA (middle), MuSiQue (right). IRCoT outperforms OneR for all model sizes, including the 0.3B model, and the difference roughly grows with model size. Note that OneR does not interleave any LM in its retrieval and hence has a fixed score.

Other Retrieval Augmented LMs

Shall We Pretrain Autoregressive Language Models with Retrieval? A Comprehensive Study

Boxin Wang^{*†1} Wei Ping^{*‡2} Peng Xu^{*2} Lawrence McAfee²
Zihan Liu² Mohammad Shoeybi² Yi Dong² Oleksii Kuchaiev²
Bo Li¹ Chaowei Xiao^{2,3} Anima Anandkumar² Bryan Catanzaro²

| Model Name | #/ Retrieval Tokens | When to Involve Retrieval | Architecture | Initialization | Re-indexing |
|----------------------------|---------------------|---------------------------|-----------------|-------------------------------|-------------|
| RETRO (Borgeaud et al.) | $O(10^{12})$ | Pretraining | decoder-only | From Scratch / Pretrained GPT | No |
| Atlas (Izacard et al.) | $O(10^9)$ | Pretraining | encoder-decoder | Pretrained T5 | Yes |
| REALM (Guu et al.) | $O(10^9)$ | Pretraining | encoder-only | Pretrained BERT | Yes |
| RAG (Lewis et al.) | $O(10^9)$ | Fine-tuning | encoder-decoder | Pretrained BART | No |
| DPR (Karpukhin et al.) | $O(10^9)$ | Fine-tuning | encoder-only | Pretrained BERT | No |
| FiD (Izacard and Grave) | $O(10^9)$ | Fine-tuning | encoder-decoder | Pretrained T5 | No |
| KNN-LM (Khandelwal et al.) | $O(10^9)$ | Inference | decoder-only | Pretrained GPT | No |

Table 1: Comparison of different retrieval-augmented models in terms of #/ retrieval tokens, which stage to incorporate retrieval into LMs, the architecture of the backbone LM, whether it requires initialization from the existing LM checkpoint, and whether it requires expensive re-indexing. RETRO is the most scalable retrieval-augmented LM due to its chunk-level retrieval and scalable decoder-only autoregressive LM backbone (Thoppilan et al., 2022; Brown et al., 2020; Smith et al., 2022; Chowdhery et al., 2022) without expensive retrieval index refresh.

Takeaways



- CoT (et al) can be viewed as general context expansion technique used in conjunction with other methods discussed
- Code conditioning provides useful signal
- Use tools like internet search / wikipedia search and combine with CoT techniques
 - Inference v/s training time

Ethical Considerations

On Second Thought, Let's Not Think Step by Step! Bias and Toxicity in Zero-Shot Reasoning

Omar Shaikh[▲], Hongxin Zhang[✉], William Held[✉], Michael Bernstein[▲], Diyi Yang[▲]
Stanford University, Shanghai Jiao Tong University, Georgia Institute of Technology
oshaikh@stanford.edu, icefox@sjtu.edu.cn, wheld3@gatech.edu
{mbernst, diyiy}@cs.stanford.edu

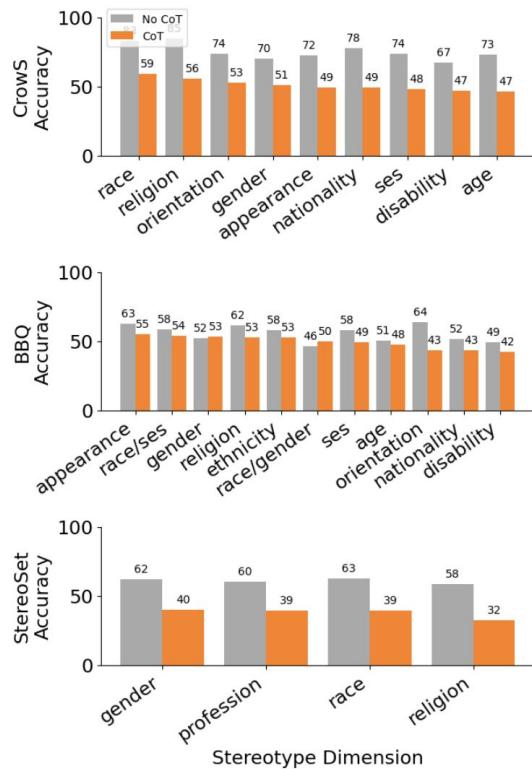


Figure 2: **Accuracy Degradations Across Dimension** for benchmark categories when using text-davinci-002. Percentages closer to 100 are better. Categories are sorted by CoT accuracy.

:)

The Cambrian explosion of LM++'s (Language models plus search, retrieval, chain of thought, etc.) will take years to *fully* settle but it'll be much clearer what works and which of the zillion products matter by the end of the year.

Until then, no one knows what's going on.

8

17

126

40.1K



Thank you!

For any questions, please reach out at gargib2@illinois.edu