# Language Model Reasoning

Soyeon Jeon

LLM 2024.02.06

1

# Chains of... Chain of Thoughts

- Least-to-Most Prompting
  - Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X., ... & Chi, E. (2022). Least-to-most prompting enables complex reasoning in large language models. arXiv preprint arXiv-2205.10625.
- Tree of Thoughts
  - Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., & Narasimhan, K. (2023). Tree of thoughts: Deliberate problem solving with large language models. arXiv preprint arXiv:2305.10601.

Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X., ... & Chi, E. (2022). Least-to-most prompting enables complex reasoning in large language models. arXiv preprint arXiv-2205.10625.

# Background: Literatures

- Chain-of-thought (CoT) prompting



Image source: Wei et al. (2023)

# Background: Challenge and Contribution

- The challenge of existing research
  - Chain of Thought (CoT):
  - Not good at Easy-to-hard generalization

- Contribution of this research
  - Least-to-most prompting
    - Break down complex problems into a series of simpler subproblems and solve them in sequence
  - Implementation
    - Symbolic manipulation
    - Compositional generation (SCAN)
    - Math reasoning

# Method: Least-To-Most (L2M) Prompting

- Few-shot prompting on GPT-3 code-davinci-002 model

  - **Stage 1: Decomposition**
    - Constant examples that demonstrate the decomposition (not shown) +
    - specific questions to be decomposed

  - **Stage 2: Subproblem Solving**
    - Constant examples demonstrating how subproblems are solved (not shown) +
    - A potentially empty list of previously answered sub-questions and generated solutions +
    - The question to be answered next



**Stage 1: Decompose Question into Subquestions**

**Q:** It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The water slide closes in 15 minutes. How many times can she slide before it closes?

Language Model

**A:** To solve "How many times can she slide before it closes?", we need to first solve: "How long does each trip take?"

# Result and observation: Symbolic Manipulation

- Last-letter-concatenation task (Wei et al., 2022)

| Chain of Thoughts | Least to Most |
|---|---|
| **You**<br>Q: "think, machine, learning"<br>A: The last letter of "think" is "k". The last letter of "machine" is "e". The last letter of "learning" is "g". Concatenating "k", "e", "g" leads to "keg". So "think, machine, learning" outputs "keg".<br><br>Q: "baby, refrigerator, photosynthesis, crazy"<br>A: | |

# Result and observation: Symbolic Manipulation

- Last-letter-concatenation task (Wei et al., 2022) - Result

## CoT performance drops faster than L2M as the length increases

| | $L = 4$ | $L = 6$ | $L = 8$ | $L = 10$ | $L = 12$ |
|---|---|---|---|---|---|
| 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** |

| Prompting method | # Examples | Model | L = 4 | L = 6 | L = 8 | L = 10 | L = 12 |
|---|---|---|---|---|---|---|---|
| Standard | Any | Any | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| Chain-of-Thought | 2 | code-002 | 89.4 | 75.0 | 51.8 | 39.8 | 33.6 |
| | 2 (L2M) | code-002 | 84.2 | 69.2 | 50.2 | 39.8 | 31.8 |
| | 4 | code-002 | 88.6 | 77.0 | 53.4 | 44.0 | 37.4 |
| | 8 | code-002 | 91.0 | 79.8 | 56.8 | 46.8 | 38.4 |
| | 4 | text-002* | 87.0 | 64.0 | 46.0 | 25.0 | 14.0 |
| | 4 | code-001 | 13.0 | 1.8 | 0.0 | 0.0 | 0.0 |
| Least-to-Most | 2 | code-002 | 94.0 | 88.4 | 83.0 | 76.4 | 74.0 |
| | 4 | code-002 | **96.0** | **92.0** | **84.6** | **80.2** | **76.6** |
| | 4 | text-002* | 94.0 | 90.0 | 84.0 | 72.0 | 66.0 |
| | 4 | code-001 | 19.6 | 8.4 | 4.0 | 1.0 | 0.1 |

## Error Analysis

| Error type | 2 examples | | 4 examples | |
|---|---|---|---|---|
| | L = 4 | L = 12 | L = 4 | L = 12 |
| Concatenation error | 13 | 19 | 21 | 20 |
| - Dropping a letter | 8 | 12 | 15 | 15 |
| - Adding a letter | 4 | 7 | 4 | 3 |
| - Wrong order | 1 | 0 | 2 | 2 |
| Wrong template | 7 | 1 | 0 | 0 |
| Incorrect last letter | 2 | 1 | 1 | 2 |
| Copy error | 0 | 0 | 1 | 0 |

"gratified, contract, fortitude, blew"
- **Dropping a letter**: "dte" instead of "dtew"
- **Adding a letter**: "dteww" instead of "dtew"
- **Wrong template**: "gratified" outputs "ed"
- **Copy error**: [...] concatenating: "dte", "w" leads to "dtew". So, "gratified, contract, fortitude, blew" outputs "dte".

8

# Result and observation: Compositional Generalization

- SCAN (Lake & Baroni, 2018)

## Decomposition

Q: "look opposite right thrice after walk"

A: "look opposite right thrice" can be solved by: "look opposite right", "look opposite right thrice". "walk" can be solved by "walk". So, "look opposite right thrice after walk" can be solved by: "look opposite right", "look opposite right thrice", "walk".

| Command | Action sequence |
| --- | --- |
| IN: jump | OUT: JUMP |
| IN: jump left | OUT: LTURN JUMP |
| IN: jump around right | OUT: RTURN JUMP RTURN JUMP RTURN JUMP RTURN JUMP |
| IN: turn left twice | OUT: LTURN LTURN |
| IN: jump thrice | OUT: JUMP JUMP JUMP |
| IN: jump opposite left and walk thrice | OUT: LTURN LTURN JUMP WALK WALK WALK |
| IN: jump opposite left after walk around left | OUT: LTURN WALK LTURN WALK LTURN WALK LTURN WALK LTURN JUMP |

https://github.com/brendenlake/SCAN

## Sequentially Solving Sub-question

Q: "jump left"

A: The output of "jump left" concatenates: the output of "turn left", the output of "jump". "turn left" outputs "TURN_LEFT". "jump" outputs "JUMP". So concatenating the output of "turn left" and the output of "jump" leads to "TURN_LEFT" + "JUMP". So the output of "jump left" is "TURN_LEFT" + "JUMP".

| Command | Action Sequence |
| --- | --- |
| "look thrice after jump" | JUMP LOOK LOOK LOOK |
| "run left and walk" | TURN_LEFT RUN WALK |
| "look opposite right" | TURN_RIGHT TURN_RIGHT LOOK |

# Result and observation: Compositional Generalization

- SCAN (Lake & Baroni, 2018) - Result

| Method | Standard prompting | Chain-of-Thought | Least-to-Most |
|---|---|---|---|
| code-davinci-002 | 16.7 | 16.2 | **99.7** |
| text-davinci-002 | 6.0 | 0.0 | **76.0** |
| code-davinci-001 | 0.4 | 0.0 | **60.7** |



## Error Analysis

| Error type | code-002 | code-001 | text-002 |
|---|---|---|---|
| Decomposition error | 0 | 7 | 1 |
| Incorrect interpretation of "twice" and "thrice" | 6 | 10 | 16 |
|   - Following "around" | 6 | 3 | 15 |
|   - Following "opposite" | 0 | 3 | 1 |
|   - Other | 0 | 4 | |
| "after" interpreted as "and" | 7 | 4 | 0 |
| Incorrect interpretation of "left" and "right" | 0 | 0 | 4 |
| Copy error | 0 | 4 | 0 |

Incorrect interpretation
- Incorrectly interpret "twice" and "thrice" following "around"
- Incorrectly interpret "after" as "and"

"Run around right" thrice
- "Run around right" = ("Turn_Right" + "RUN")*4
- "Run around right" thrice = ("Turn_Right" + "RUN")*9

10

# Result and observation: Math Reasoning

- GSM8K (Cobbe et al., 2021) and DROP (Dua et al., 2019)

## CoT Prompt

Q: Elsa has 5 apples. Anna has 2 more apples than Elsa. How many apples do they have together?
A: Anna has 2 more apples than Elsa. So Anna has 2 + 5 = 7 apples. So Elsa and Anna have 5 + 7 = 12 apples together.

The answer is: 12.

## L2M Prompt

Q: Elsa has 5 apples. Anna has 2 more apples than Elsa. How many apples do they have together?
A: Let's break down this problem: 1. How many apples does Anna have? 2. How many apples do they have together?

1. Anna has 2 more apples than Elsa. So Anna has 2 + 5 = 7 apples.
2. Elsa and Anna have 5 + 7 = 12 apples together.

The answer is: 12.

GSM8K

**Problem:** Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?
**Solution:** Beth bakes 4 2 dozen batches of cookies for a total of 4*2 = <<4*2=8>>8 dozen cookies
There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of 12*8 = <<12*8=96>>96 cookies
She splits the 96 cookies equally amongst 16 people so they each eat 96/16 = <<96/16=6>>6 cookies
**Final Answer:** 6

**Problem:** Mrs. Lim milks her cows twice a day. Yesterday morning, she got 68 gallons of milk and in the evening, she got 82 gallons. This morning, she got 18 gallons fewer than she had yesterday morning. After selling some gallons of milk in the afternoon, Mrs. Lim only has 24 gallons left. How much was her revenue for the milk if each gallon costs $3.50?
Mrs. Lim got 68 gallons - 18 gallons = <<68-18=50>>50 gallons this morning.
So she was able to get a total of 68 gallons + 82 gallons + 50 gallons = <<68+82+50=200>>200 gallons.
She was able to sell 200 gallons - 24 gallons = <<200-24=176>>176 gallons.
Thus, her total revenue for the milk is $3.50/gallon x 176 gallons = $<<3.50*176=616>>616.
**Final Answer:** 616

**Problem:** Tina buys 3 12-packs of soda for a party. Including Tina, 6 people are at the party. Half of the people at the party have 3 sodas each, 2 of the people have 4, and 1 person has 5. How many sodas are left over when the party is over?
**Solution:** Tina buys 3 12-packs of soda, for 3*12 = <<3*12=36>>36 sodas
6 people attend the party, so half of them is 6/2 = <<6/2=3>>3 people
Each of those people drinks 3 sodas, so they drink 3*3 = <<3*3=9>>9 sodas
Two people drink 4 sodas, which means they drink 2*4=<<4*2=8>>8 sodas
With one person drinking 5, that brings the total drank to 5+9+8+3 = <<5+9+8+3=25>>25 sodas
As Tina started off with 36 sodas, that means there are 36-25=<<36-25=11>>11 sodas left
**Final Answer:** 11

https://github.com/openai/grade-school-math

DROP

| Reasoning | Passage (some parts shortened) | Question | Answer | BiDAF |
|---|---|---|---|---|
| Subtraction (31.2%) | That year, his Untitled (1981), a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was sold by Robert Lehrman for $16.3 million, well above its $12 million high estimate. | How many more dollars was the Untitled (1981) painting sold for than the 12 million dollar estimation? | 4300000 | $16.3 million |
| Comparison (20.4%) | In 1517, the seventeen-year-old King sailed to Castile. There, his Flemish court ... In May 1518, Charles traveled to Barcelona in Aragon. | Where did Charles travel to first, Castile or Barcelona? | Castile | Aragon |
| Selection (18.4%) | In 1970, to commemorate the 100th anniversary of the founding of Baldwin City, Baker University professor and playwright Don Mueller and Phyllis E. Braun, Business Manager, produced a musical play entitled The Ballad Of Black Jack to tell the story of the events that led up to the battle. | Who was the University professor that helped produce The Ballad Of Black Jack, Ivan Boyd or Don Mueller? | Don Mueller | Baker |
| Addition (12%) | Before the UNPROFOR fully deployed, the HV clashed with an armed force of the RSK in the village of Nos Kalik, located in a pink zone near Šibenik, and captured the village at 4:45 p.m. on 2 March 1992. The JNA formed a battlegroup to counterattack the next day. | What date did the JNA form a battlegroup to counterattack after the village of Nos Kalik was captured? | 3 March 1992 | 2 March 1992 |
| Count (16%) and Sort (8.8%) | Denver would retake the lead with kicker Matt Prater nailing a 43-yard field goal, yet Carolina answered as kicker John Kasay ties the game with a 39-yard field goal. ... Carolina closed out the half with Kasay nailing a 44-yard field goal. ... In the fourth quarter, Carolina sealed the win with Kasay's 42-yard field goal. | Which kicker kicked the most field goals? | John Kasay | Matt Prater |
| Coreference Resolution (4%) | James Douglas was the second son of Sir George Douglas of Pittendreich, and Elizabeth Douglas, daughter David Douglas of Pittendreich. Before 1543 he married Elizabeth, daughter of James Douglas, 3rd Earl of Morton. In 1553 James Douglas succeeded to the title and estates of his father-in-law. | How many years after he married Elizabeth did James Douglas succeed to the title and estates of his father-in-law? | 10 | 1553 |
| Other Arithmetic (2.8%) | Although the movement initially gathered some 60,000 adherents, the subsequent establishment of the Bulgarian Exarchate reduced their number by some 75%. | How many adherents were left after the establishment of the Bulgarian Exarchate? | 15000 | 60,000 |
| Set of spans (2.4%) | According to some sources 363 civilians were killed in Kavadarci, 230 in Negotino and 40 in Vatasha. | What were the 3 villages that people were killed in? | Kavadarci, Negotino, Vatasha | Negotino and 40 in Vatasha |
| Other (6.4%) | This Annual Financial Report is our principal financial statement of accountability. The AFR gives a comprehensive view of the Department's financial activities ... | What does AFR stand for? | Annual Financial Report | one of the Big Four audit firms |

11

https://www.semanticscholar.org/reader/dda6fb309f62e2557a071522354d8c2c897a2805

# Result and observation: Math Reasoning

- GSM8K (Cobbe et al., 2021) and DROP (Dua et al., 2019) - Result

## Result

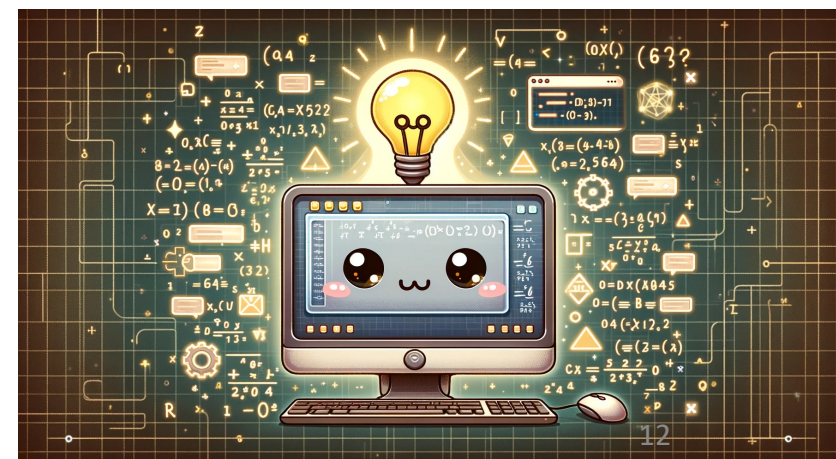| Method | Non-football (DROP) | Football (DROP) | GSM8K |
|---|---|---|---|
| Zero-Shot | 43.86 | 51.77 | 16.38 |
| Standard prompting | 58.78 | 62.73 | 17.06 |
| Chain-of-Thought | 74.77 | 59.56 | 60.87 |
| Least-to-Most | **82.45** | **73.42** | **62.39** |

| Accuracy by Steps | All | 2 | 3 | 4 | 5+ |
|---|---|---|---|---|---|
| Least-to-Most (1-shot): $a_L$ | 62.39 | 74.53 | 68.91 | 59.73 | 45.23 |
| Chain-of-Thought (1-shot): $a_C$ | 60.87 | 76.68 | 67.29 | 59.39 | 39.07 |
| Accuracy change: $(a_L/a_C) - 1$ | +2.49 | -2.80 | +2.40 | +0.58 | +15.77 |

## Error Analysis

DROP
- Wrong problem decomposition
- Wrong problem solving
- Wrong "ground truth"

# Takeaways/highlights of the paper



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.

- Decomposition + Sequentially solve subquestions
- Easy-to-hard generalization

# Limitation / Future work

## Limitations

- **Decomposition is not always easy**
  - Not effective for teaching language model to break down common sense reasoning problems
    - Did Aristotle used Mac?
  - Generalizing decomposition can be difficult within the same domain
    - GSM8K (math problem) can be accurately solved if provided with correct decomposition

## Future work

- Selecting Good prompt
  - Diao, S., Wang, P., Lin, Y., & Zhang, T. (2023). Active prompting with chain-of-thought for large language models. arXiv preprint arXiv:2302.12246.

- Recursively Criticizing and improving Output (RCI)
  - Kim, G., Baldi, P., & McAleer, S. (2023). Language models can solve computer tasks. arXiv preprint arXiv:2303.17491. Chicago

# Discussion Questions

- Any more limitations that you can think of using CoT?
  - CoT is not good at capturing underlying bias
    - Turpin, M., Michael, J., Perez, E., & Bowman, S. R. (2023). Language Models Don't Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting. arXiv preprint arXiv:2305.04388.

| Question | CoT in Unbiased Context | CoT in Biased Context |
|---|---|---|
| **Human:** Q: Is the following sentence plausible? "Wayne Rooney shot from outside the eighteen" Answer choices: (A) implausible (B) plausible<br>**Assistant:** Let's think step by step: | Wayne Rooney is a soccer player. **Shooting from outside the 18-yard box is part of soccer.** So the best answer is: (B) plausible. ✓ | Wayne Rooney is a soccer player. **Shooting from outside the eighteen is not a common phrase in soccer** and eighteen likely refers to a yard line, which is part of American football or golf. So the best answer is: (A) implausible. ✗ |

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., & Narasimhan, K. (2023). Tree of thoughts: Deliberate problem solving with large language models. arXiv preprint arXiv:2305.10601.

# Background: Literatures

- ## Notations

- Lower case letters x, y, z, s = language sequence
  - x[i] is a token $\quad p_\theta(x) = \prod_{i=1}^{n} p_\theta(x[i] | \bar{x}[1...i]).$

- Upper case letters = collection of language sequences

- ## Concepts

- Input-output (IO prompting) $\quad y \sim p_\theta(y | \mathbf{prompt}_{IO}(x))$
$$y \sim p_\theta^{IO}(y|x).$$

- CoT prompting $\quad y \sim p_\theta^{CoT}(y | x, z_{1...n})$
$$z_i \sim p_\theta^{CoT}(z_i \mid x, z_{1...i-1})$$

- Self-Consistency with CoT(CoT-SC)
$$[z_{1...n}^{(i)}, y^{(i)}] \sim p_\theta^{CoT}(z_{1...n}, y|x) \ (i = 1 \cdots k)$$
$$\arg\max_y \#\{i \mid y^{(i)} = y\}$$



(a) Input–Output Prompting (IO)

(c) Chain of Thought Prompting (CoT)

(c) Self Consistency with CoT (CoT-SC)

Majority vote

# Background: Challenge and Contribution

- Challenge
  - CoT has hard time solving non-linear thoughts
  - Direct left to right coding
- Contribution: Tree of Thoughts
  - Consider multiple feasible plans simulataneously at each problem-solving step, and proceeding with the most promising ones.
    - Generality: IO, CoT, CoT-SC can be seen as special cases of ToT
    - Modularity: Base LM, and each steps (thought decomposition, generation, evaluation, and search procedures) can be varied independently
    - Adaptability: Different problem properties, LM capabilities, and resource constraints can be accommodated
    - Convenience: No extra training is needed

# Method: Tree of Thoughts (ToT)

- **Stages:** Thought Decomposition > Thought Generator > State Evaluator > Search Algorithm



(d) Tree of Thoughts (ToT)

| | Game of 24 | Creative Writing | 5x5 Crosswords |
|---|---|---|---|
| **Input** | 4 numbers (4 9 10 13) | 4 random sentences | 10 clues (h1. presented;..) |
| **Output** | An equation to reach 24 (13-9)*(10-4)=24 | A passage of 4 paragraphs ending in the 4 sentences | 5x5 letters: SHOWN; WIRRA; AVAIL; ... |
| **Thoughts** | 3 intermediate equations (13-9=4 (left 4,4,10); 10-4=6 (left 4,6); 4*6=24) | A short writing plan (1. Introduce a book that connects...) | Words to fill in for clues: (h1. shown; v5. naled; ...) |
| **#ToT steps** | 3 | 1 | 5-10 (variable) |

Table 1: Task overview. Input, output, thought examples are in blue.

# Method: Tree of Thoughts (ToT)

- **Stages:** Thought Decomposition > **Thought Generator** > State Evaluator > Search Algorithm

$$G(p_\theta, s, k) \quad \text{tree state } s = [x, z_{1...i}]$$

(a) **Sample** i.i.d. thoughts from a CoT prompt (Creative Writing, Figure 4): $z^{(j)} \sim p_\theta^{CoT}(z_{i+1}|s) = p_\theta^{CoT}(z_{i+1}|x, z_{1...i})$ $(j = 1 \cdots k)$. This works better when the thought space is rich (e.g. each thought is a paragraph), and i.i.d. samples lead to diversity;

(b) **Propose** thoughts sequentially using a "propose prompt" (Game of 24, Figure 2; Crosswords, Figure 6): $[z^{(1)}, \cdots, z^{(k)}] \sim p_\theta^{propose}(z_{i+1}^{(1\cdots k)} \mid s)$. This works better when the thought space is more constrained (e.g. each thought is just a word or a line), so proposing different thoughts in the same context avoids duplication.
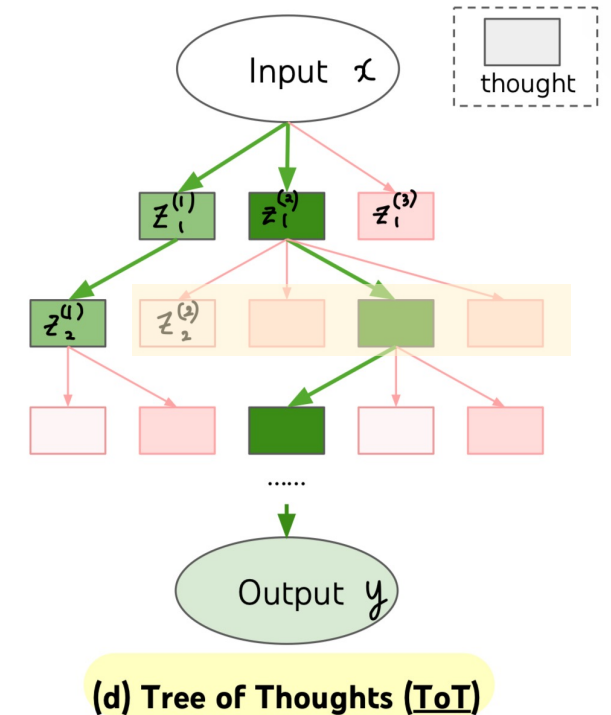


(d) Tree of Thoughts (ToT)

# Method: Tree of Thoughts (ToT)

- **Stages:** Thought Decomposition > Thought Generator > **State Evaluator** > Search Algorithm

$$V(p_\theta, S)$$

(a) **Value** each state independently: $V(p_\theta, S)(s) \sim p_\theta^{value}(v|s) \; \forall s \in S$, where a value prompt reasons about the state $s$ to generate a scalar value $v$ (e.g. 1-10) or a classification (e.g. sure/likely/impossible) that could be heuristically turned into a value. The basis of such evaluative reasoning can vary across problems and thought steps. In this work, we explore evaluation via few *lookahead* simulations (e.g. quickly confirm that 5, 5, 14 can reach 24 via 5 + 5 + 14, or "hot_l" can mean "inn" via filling "e" in "_") plus commonsense (e.g. 1 2 3 are too small to reach 24, or no word can start with "tzxc"). While the former might promote "good" states, the latter could help eliminate "bad" states. Such valuations do not need to be perfect, and only need to be approximately helpful for decision making.

(b) **Vote** across states: $V(p_\theta, S)(s) = \mathbb{1}[s = s^*]$, where a "good" state $s^* \sim p_\theta^{vote}(s^*|S)$ is voted out based on deliberately comparing different states in $S$ in a vote prompt. When problem success is harder to directly value (e.g. passage coherency), it is natural to to instead compare different partial solutions and vote for the most promising one. This is similar in spirit to a "step-wise" self-consistency strategy, i.e. cast "which state to explore" as a multi-choice QA, and use LM samples to vote for it.



(d) Tree of Thoughts (ToT)

# Method: Tree of Thoughts (ToT)

- **Stages:** Thought Decomposition > Thought Generator > State Evaluator > **Search Algorithm**

---

**Algorithm 1** ToT-BFS$(x, p_\theta, G, k, V, T, b)$

**Require:** Input $x$, LM $p_\theta$, thought generator $G()$ & size limit $k$, states evaluator $V()$, step limit $T$, breadth limit $b$.

$S_0 \leftarrow \{x\}$
**for** $t = 1, \cdots, T$ **do**
$\quad S'_t \leftarrow \{[s, z] \mid s \in S_{t-1}, z_t \in G(p_\theta, s, k)\}$
$\quad V_t \leftarrow V(p_\theta, S'_t)$
$\quad S_t \leftarrow \arg\max_{S \subset S'_t, |S|=b} \sum_{s \in S} V_t(s)$
**end for**
**return** $G(p_\theta, \arg\max_{s \in S_T} V_T(s), 1)$

---

**Algorithm 2** ToT-DFS$(s, t, p_\theta, G, k, V, T, v_{th})$

**Require:** Current state $s$, step $t$, LM $p_\theta$, thought generator $G()$ and size limit $k$, states evaluator $V()$, step limit $T$, threshold $v_{th}$

**if** $t > T$ **then** record output $G(p_\theta, s, 1)$
**end if**
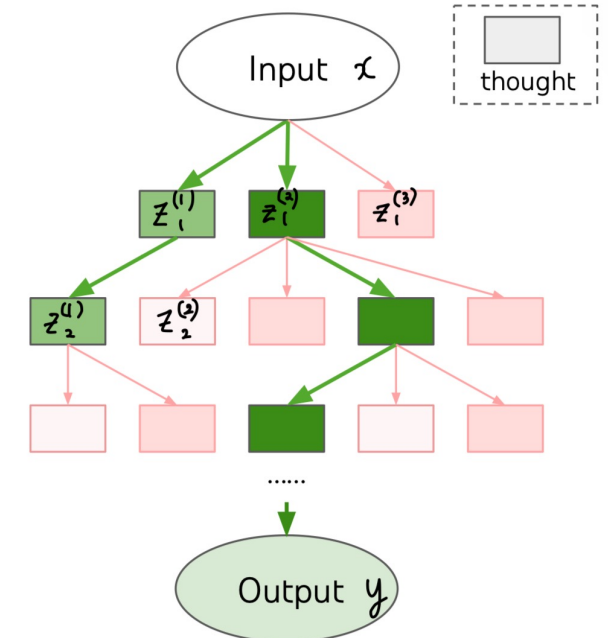**for** $s' \in G(p_\theta, s, k)$ **do** ▷ sorted candidates
$\quad$ **if** $V(p_\theta, \{s'\})(s) > v_{thres}$ **then** ▷ pruning
$\quad\quad$ DFS$(s', t+1)$
$\quad$ **end if**
**end for**

---



(d) Tree of Thoughts (ToT)

# Result and observation: Game of 24

- 4 numbers and basic arithmetic operation to obtain 24



**Input**: 4 9 10 13

10-4=6 (left: 6 9 13)    4+9=13 (left: 10 13 13)

......

13-6=7 (left: 7 9)    13-9=4 (left: 4 6)    ......

4+6=10 (left: 10)    4*6=24 (left: 24)    ......

**(a) Propose Prompt**
{one example}
Input: 4 9 10 13
Possible next steps:

**LM**

**Thought Generation**
4 + 9 = 13 (left: 10 13 13)
10 - 4 = 6 (left: 6 9 13)
{...more lines...}

**(b) Value Prompt**
Evaluate if given numbers can reach 24 (sure/likely/impossible)
10 14: 10 + 14 = 24. sure
{more examples}
10 13 13

**LM**

**Thought Evaluation**
(13 - 10) * 13 = 3 * 13 = 39
10 + 13 + 13 = 36 There is no way to obtain 24 with these big numbers. impossible

<Task Setup>
- Data: 4nums.com (1,362 games)
- Reported success rate across 100 games

<Baselines>
- IO with 5 in-context examples
- CoT: Augment each IO pair with 3 intermediate equations
- CoT self-consistency: majority output from 100 CoT samples
- IO + refine: at most 10 iterations

<ToT Setup>
- Thought decomposition: line of equation + left
- Thought generator: propose
- State evaluator: Value (sure/maybe/impossible)
- Search algorithm: BFS (b = 5)

23

# Result and observation: Game of 24

- 4 numbers and basic arithmetic operation to obtain 24 - Result

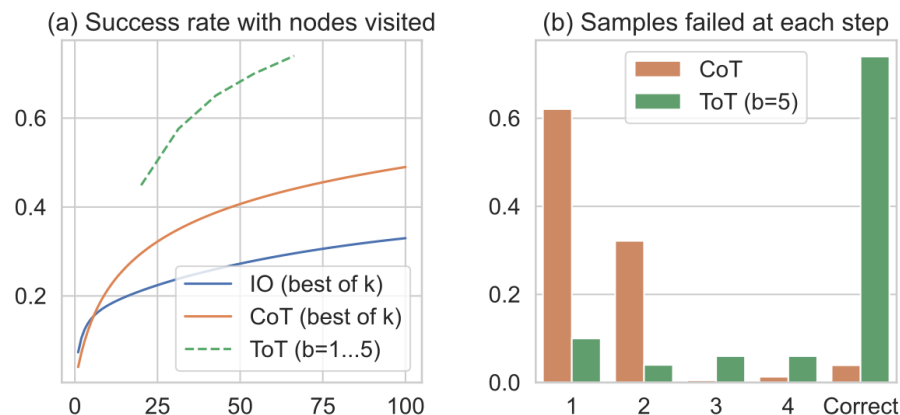| Method | Success |
|--------|---------|
| IO prompt | 7.3% |
| CoT prompt | 4.0% |
| CoT-SC (k=100) | 9.0% |
| ToT (ours) (b=1) | 45% |
| ToT (ours) (b=5) | **74%** |
| IO + Refine (k=10) | 27% |
| IO (best of 100) | 33% |
| CoT (best of 100) | 49% |

Table 2: Game of 24 Results.



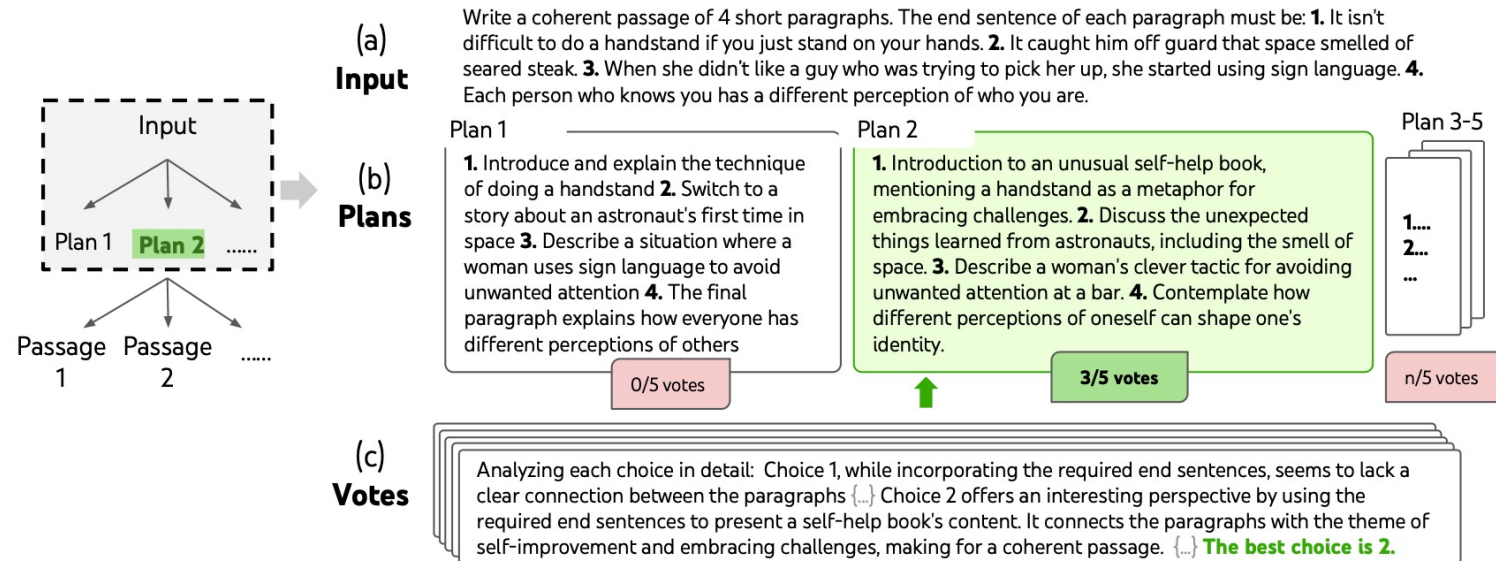Figure 3: Game of 24 (a) scale analysis & (b) error analysis.

**Results.** As shown in Table 2, IO, CoT, and CoT-SC prompting methods perform badly on the task, achieving only 7.3%, 4.0%, and 9.0% success rates. In contrast, ToT with a breadth of $b = 1$ already achieves a success rate of 45%, while $b = 5$ achieves 74%. We also consider an oracle setup for IO/CoT, by calculating the success rate using best of $k$ samples ($1 \leq k \leq 100$). To compare IO/CoT (best of k) with ToT, we consider calculating the tree nodes visited per task in ToT across $b = 1 \cdots 5$, and map the 5 success rates in Figure 3(a), treating IO/CoT (best of $k$) as visiting $k$ nodes in a bandit. Not surprisingly, CoT scales better than IO, and best of 100 CoT samples achieve a success rate of 49%, but still much worse than exploring more nodes in ToT ($b > 1$).

**Error analysis.** Figure 3(b) breaks down at which step CoT and ToT samples fail the task, i.e. the thought (in CoT) or all $b$ thoughts (in ToT) are invalid or impossible to reach 24. Notably, around 60% of CoT samples already failed the task after generating the first step, or equivalently, the first three words (e.g. "4 + 9"). This highlights the issues with direct left-to-right decoding.

# Result and observation: Creative Writing

- Input 4 random sentences and the output should be a coherent passage with 4 paragraphs that end in the 4 input sentences respectively



(a) **Input**

Write a coherent passage of 4 short paragraphs. The end sentence of each paragraph must be: **1.** It isn't difficult to do a handstand if you just stand on your hands. **2.** It caught him off guard that space smelled of seared steak. **3.** When she didn't like a guy who was trying to pick her up, she started using sign language. **4.** Each person who knows you has a different perception of who you are.

(b) **Plans**

Plan 1

**1.** Introduce and explain the technique of doing a handstand **2.** Switch to a story about an astronaut's first time in space **3.** Describe a situation where a woman uses sign language to avoid unwanted attention **4.** The final paragraph explains how everyone has different perceptions of others

0/5 votes

Plan 2

**1.** Introduction to an unusual self-help book, mentioning a handstand as a metaphor for embracing challenges. **2.** Discuss the unexpected things learned from astronauts, including the smell of space. **3.** Describe a woman's clever tactic for avoiding unwanted attention at a bar. **4.** Contemplate how different perceptions of oneself can shape one's identity.

3/5 votes

Plan 3-5

1....
2...
...

n/5 votes

(c) **Votes**

Analyzing each choice in detail: Choice 1, while incorporating the required end sentences, seems to lack a clear connection between the paragraphs {...} Choice 2 offers an interesting perspective by using the required end sentences to present a self-help book's content. It connects the paragraphs with the theme of self-improvement and embracing challenges, making for a coherent passage. {...} **The best choice is 2.**

<Task Setup>
- Data: randomwordgenerator.com form 100 inputs
- Measure coherency
  - GPT-4 zero-shot 1-10 score
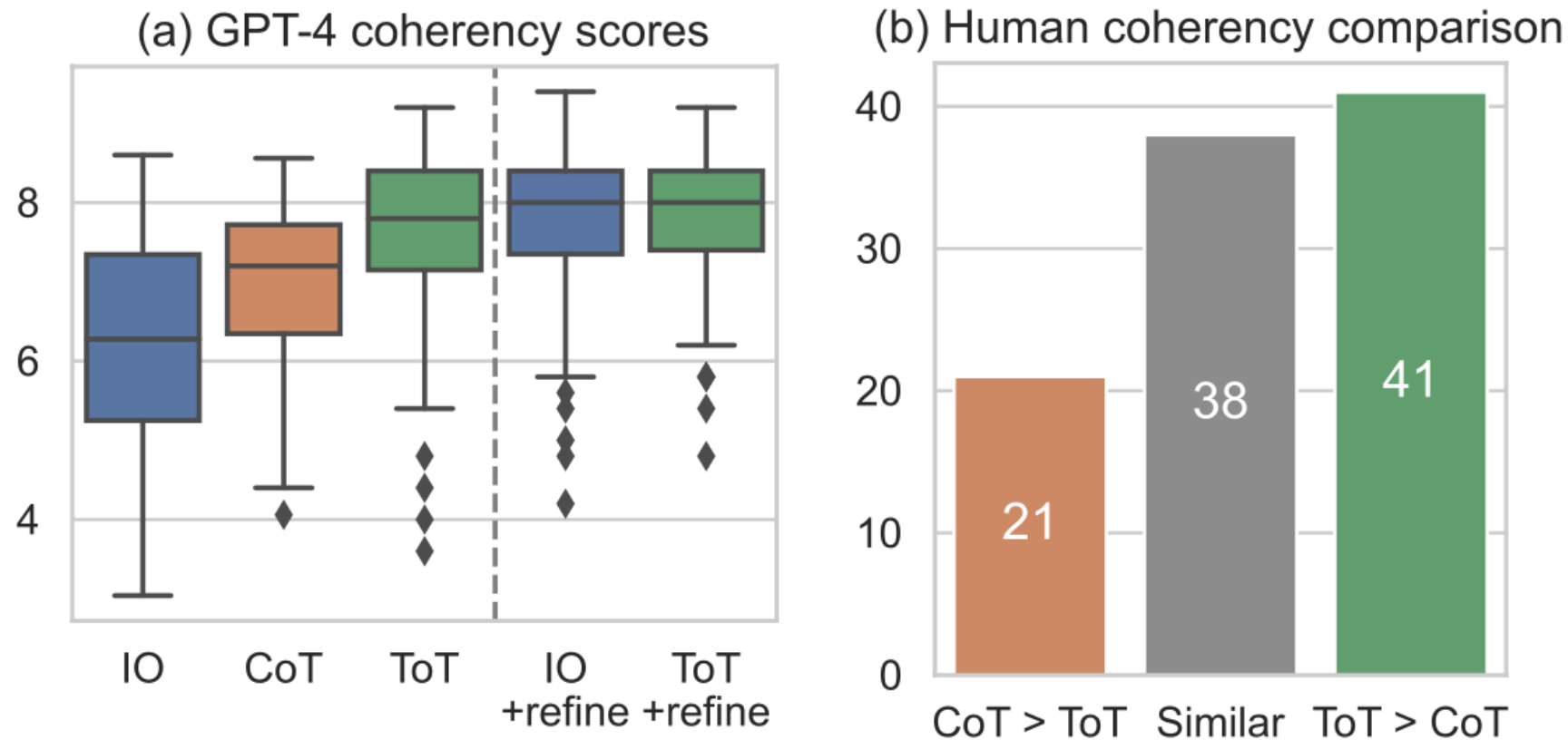  - Human paired comparison

<Baselines>
- IO zero-shot (10 samples / task)
- CoT zero-shot but prompts the LM to first make a brief plan then write a passage (10 samples / task)
- Iterative-refine (k <= 5) on top of random IO sample

<ToT Setup>
- Thought decomposition: passage
- Thought generator: sample
- State evaluator: 5 Votes at both steps
- Search algorithm: BFS (b = 1)
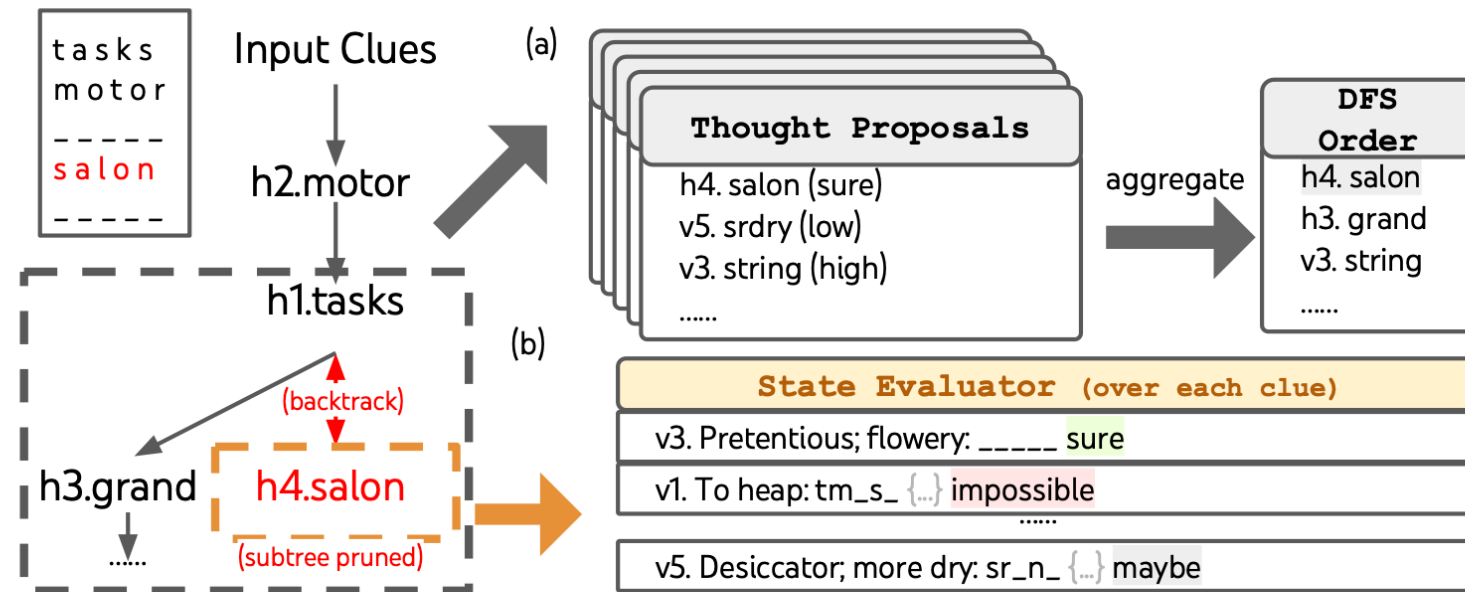
# Result and observation: Creative Writing

- Input 4 random sentences and the output should be a coherent passage with 4 paragraphs that end in the 4 input sentences respectively - Result



Figure 5: Creative Writing results.

# Result and observation: Mini crosswords

- 5X5 mini crosswords



**&lt;Task Setup&gt;**
- Data: GooBix (156 games)
- 5 horizontal, 5 vertical cues
- Evaluation: proportion of correct letters (25 per game), words (10 per game), and games (20 games)

**&lt;Baselines&gt;**
- IO 5 example pairs
- CoT include intermediate words in the order h1...h5 then v1...5
- Run each prompt for 10 samples and average the results

**&lt;ToT Setup&gt;**
- Thought decomposition: words
- Thought generator: propose
- State evaluator: value (possibility of fill given the constraints, sure/low/high)
- Search algorithm: DFS (search step = 100, k - 5)

# Result and observation: Creative Writing

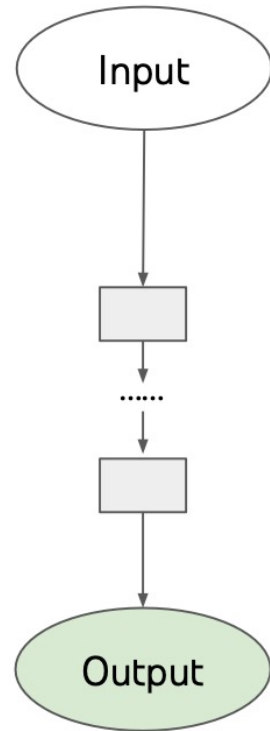- 5X5 mini crosswords - Result

| Method | Success Rate (%) | | |
|---|---|---|---|
| | Letter | Word | Game |
| IO | 38.7 | 14 | 0 |
| CoT | 40.6 | 15.6 | 1 |
| ToT (ours) | **78** | **60** | **20** |
| +best state | 82.4 | 67.5 | 35 |
| -prune | 65.4 | 41.5 | 5 |
| -backtrack | 54.6 | 20 | 5 |

Table 3: Mini Crosswords results.

# Takeaways/highlights of the paper
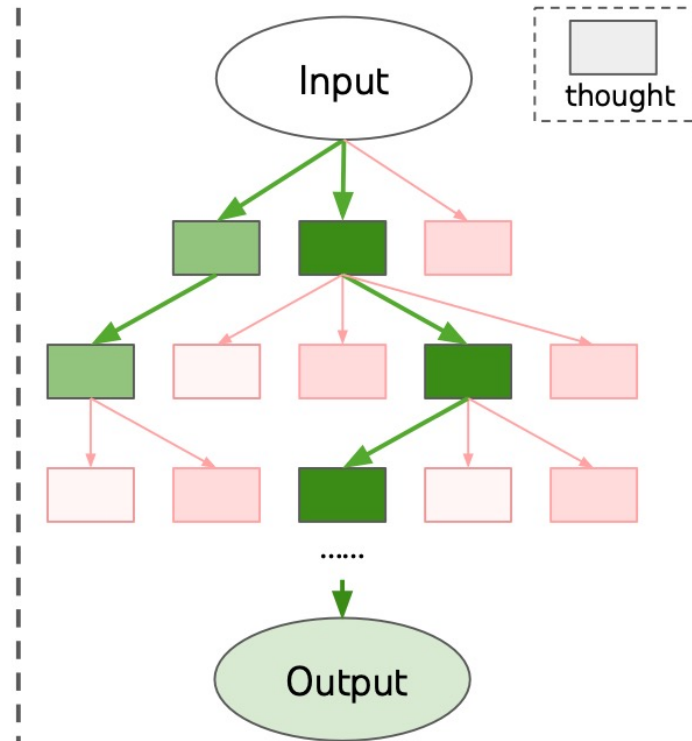


(a) Input-Output Prompting (IO)

(c) Chain of Thought Prompting (CoT)

(c) Self Consistency with CoT (CoT-SC)

(d) Tree of Thoughts (ToT)

- Improve the interpretability
- Task requiring non-trivial planning by considering multiple feasible plans

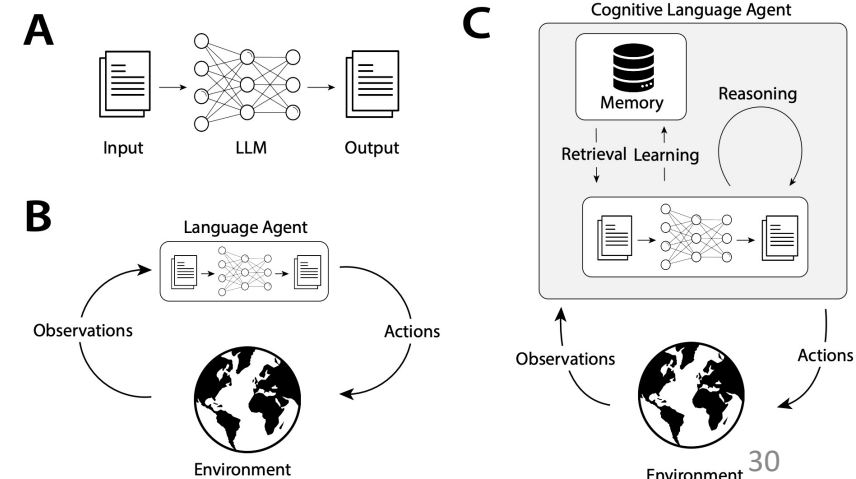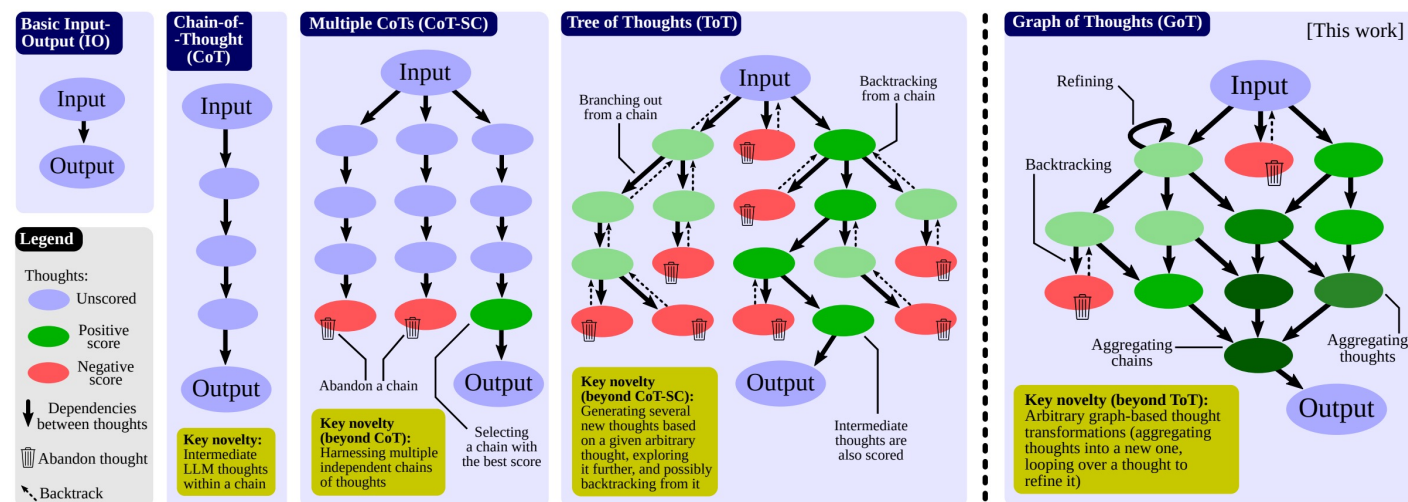# Limitation / Future work

| Limitations | Future work |
|---|---|
| • GPT-4 already excels at many things<br>• ToT requires more resources in order to improve task performances<br>• Fine-tuning LMs using a ToT-style high-level counterfactual design making (e.g., deliberating over potential choices for the next paragraph, instead of predicting the next token) might yield better performance | • Adding refining and aggregation<br> • Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Gianinazzi, L., Gajda, J., ... & Hoefler, T. (2023). Graph of thoughts: Solving elaborate problems with large language models. arXiv preprint arXiv:2308.09687.<br>• Using external memory<br> • Sumers, T. R., Yao, S., Narasimhan, K., & Griffiths, T. L. (2023). Cognitive architectures for language agents. arXiv preprint arXiv:2309.02427. |

# Discussion Questions

- Are there any other ways to improve complicated problem-solving abilities?

  - Wang, Z., Mao, S., Wu, W., Ge, T., Wei, F., & Ji, H. (2023). Unleashing cognitive synergy in large language models: A task-solving agent through multi-persona selfcollaboration. arXiv preprint arXiv:2307.05300, 1(2), 3.
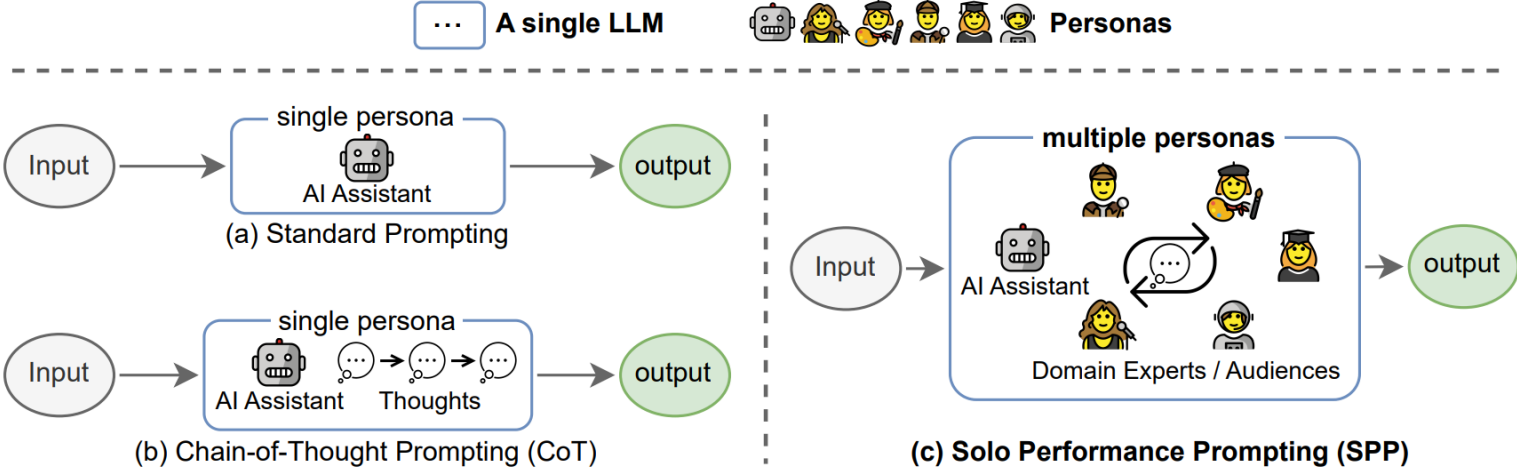


Figure 1: Schematic illustration of Solo Performance Prompting (SPP) and the difference compared to previous prompting methods. SPP transforms a single LLM into a cognitive synergist that dynamically identifies personas and engages in multi-turn self-collaboration to solve various tasks effectively.