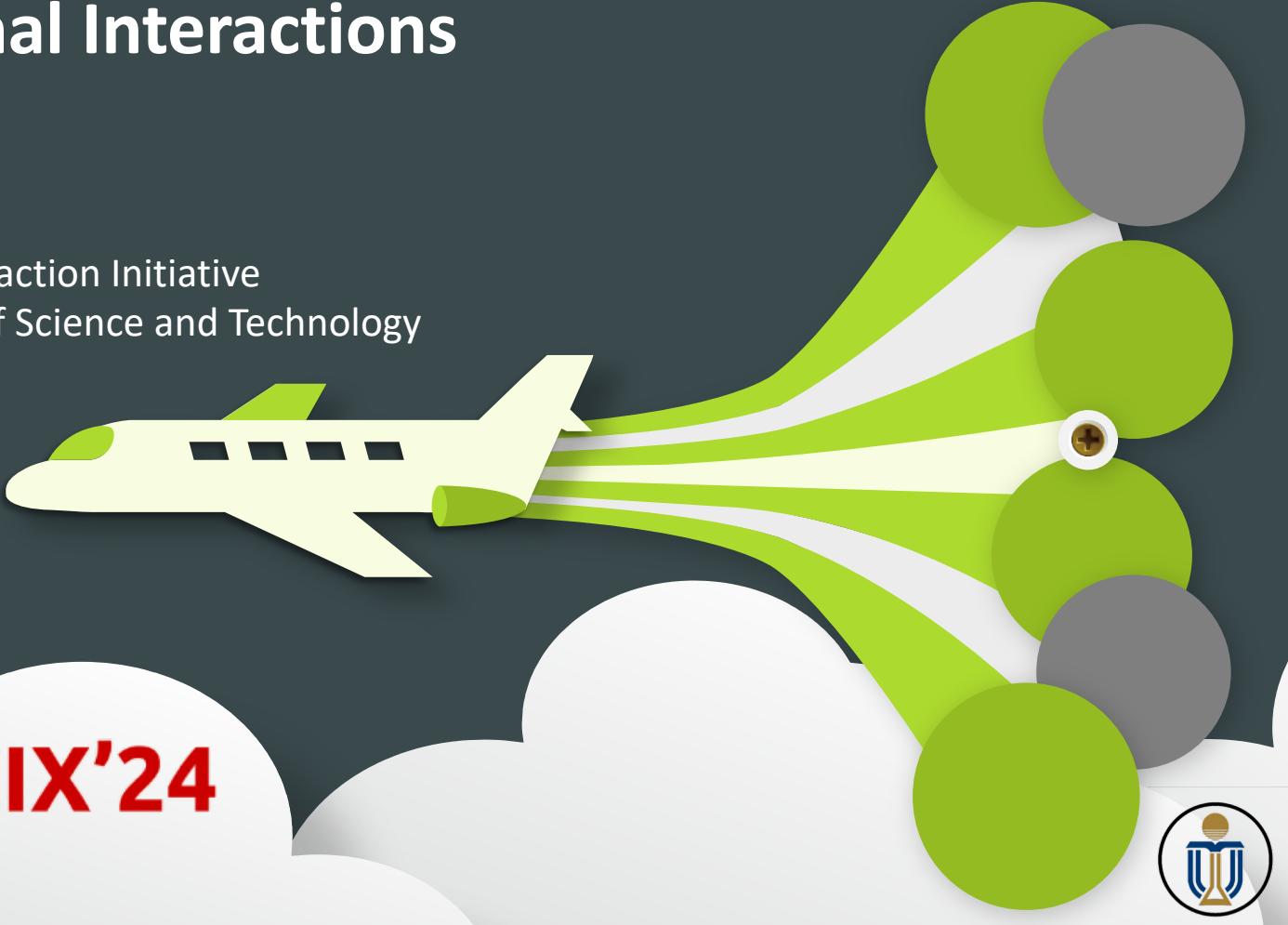


Charting the Routes of Thoughts in LLM-Empowered Conversational Interactions

Xiaojuan Ma

Human-Computer Interaction Initiative
Hong Kong University of Science and Technology
mxj@cse.ust.hk



CIX'24





Affordable and Scalable Services
Natural Interactions
Personalization



Fast Development of Foundation
Models (e.g., LLMs, VLMs, etc.)

01 Introduction

Bloomberg
Newsletter

Social Networks Want to Be Conversational AI's Killer App

Say Hi to My AI!
Introducing your fun, experimental AI sidekick exclusively on Snapchat+
ME: Can you write a haiku about my cheese-obsessed friend Lukas?
MT AI: Lukas's love for cheese, Gouda, brie, and camembert, Melts hearts, not just cheese! 🍔

Snapchat's AI chatbot Source: Snap Inc.

Adobe
Creativity & Design ▾ PDF & E-signatures ▾ Marketing & Commerce ▾ Help & Support ▾
NEWSROOM News Analyst Relations Press Contact

Adobe Brings Conversational AI to Trillions of PDFs with the New AI Assistant in Reader and Acrobat

Tuesday, February 20, 2024 09:00 AM

Adobe Brings Conversational AI to Trillions of PDFs with the New AI Assistant in Reader and Acrobat

AI Assistant in beta builds on Acrobat Liquid Mode to further unlock document intelligence with new capabilities in Reader and Acrobat

Today's release is Adobe's first step in transforming digital document experiences with generative AI for consumption and creation

Reader and Acrobat customers will have access to the full range of AI Assistant capabilities through a new add-on subscription plan when AI Assistant is out of beta

Walmart Global Tech
Home About Careers

CATHAY PACIFIC

Cathay launches advanced conversational AI powered by Fano Labs to enhance customers' digital experience

New conversational AI will enable the company's chatbots to provide more accurate responses to customers' queries

Thursday 29 September 2022 — Cathay is enhancing its customers' digital experience with the launch of an advanced conversational artificial intelligence (AI) in partnership with Fano Labs, a Hong Kong-based language AI company.

As part of its purpose to move people forward in life, Cathay is continuously pursuing the development of new technologies that enable it to give customers more choice and control over their journeys, whether on the ground, in the air or on digital channels.

Read more

Three Ways We're Using Conversational AI at Walmart

Cheryl Ainoa
SVP - New Businesses and Emerging Technology
March 3, 2023 7 min read

Customers are shifting from "do it with me" to more proactive "do it for me" experiences, and AI is a critical enabler in how we provide those experiences, whenever and wherever they need us. As a result, Walmart has been developing conversational AI capabilities for the past several years.

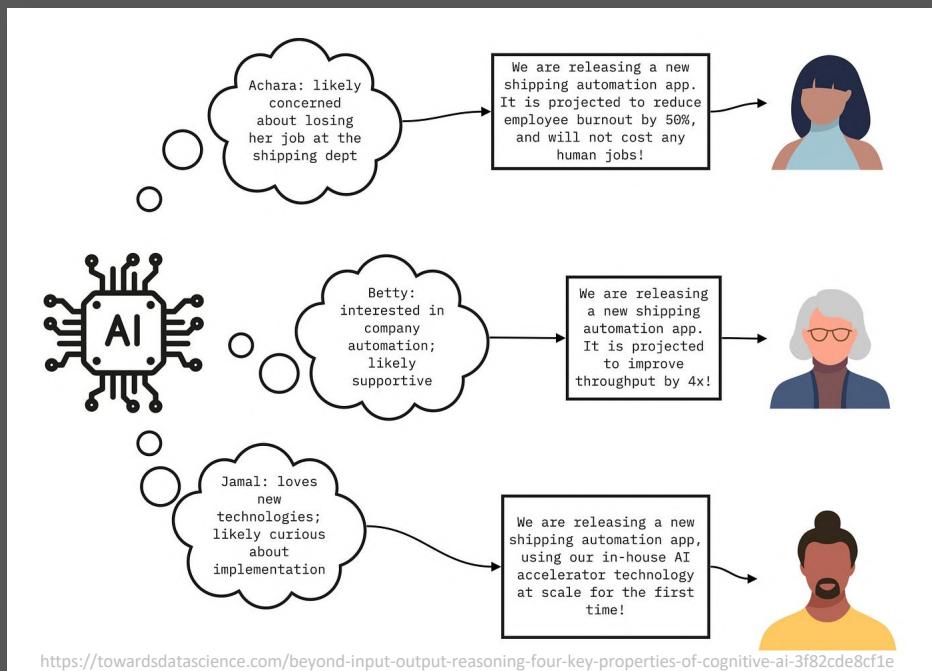
Over time, it has evolved into a robust multi-layered solution that spans across several areas of business to support our customers, members, associates, drivers and marketplace sellers, among others.

Here's a look at how we're leveraging conversational AI to help our customers and associates around the world save time and have a better experience:

Affordable and Scalable Services

Natural Interactions

Personalization



Fast Development of Foundation Models (e.g., LLMs, VLMs, etc.)

01 Introduction

Bloomberg

Newsletter

Social Networks Want to Be Conversational AI's Killer App

Say Hi to My AI!

Introducing your fun, experimental AI sidekick exclusively on Snapchat+

ME: Can you write a haiku about my cheese-obsessed friend Lukas?

MT: Lukas's love for cheese, Gouda, brie, and camembert, Melts hearts, not just cheese! 🍔

Snapchat's AI chatbot Source: Snap Inc.

Adobe

Creativity & Design ▾ PDF & E-signatures ▾ Marketing & Commerce ▾ Help & Support ▾

NEWSROOM News Analyst Relations Press Contact

Adobe Brings Conversational AI to Trillions of PDFs with the New AI Assistant in Reader and Acrobat

Tuesday, February 20, 2024 09:00 AM

Adobe Brings Conversational AI to Trillions of PDFs with the New AI Assistant in Reader and Acrobat

AI Assistant in beta builds on Acrobat Liquid Mode to further unlock document intelligence with new capabilities in Reader and Acrobat

Today's release is Adobe's first step in transforming digital document experiences with generative AI for consumption and creation

Reader and Acrobat customers will have access to the full range of AI Assistant capabilities through a new add-on subscription plan when AI Assistant is out of beta

Walmart Global Tech

Home About Careers

Three Ways We're Using Conversational AI at Walmart

Cheryl Ainoa SVP - New Businesses and Emerging Technology March 3, 2023 7 min read

Cathay launches advanced conversational AI powered by Fano Labs to enhance customers' digital experience

New conversational AI will enable the company's chatbots to provide more accurate responses to customers' queries

Thursday 29 September 2022 — Cathay is enhancing its customers' digital experience with the launch of an advanced conversational artificial intelligence (AI) in partnership with Fano Labs, a Hong Kong-based language AI company.

As part of its purpose to move people forward in life, Cathay is continuously pursuing the development of new technologies that enable it to give customers more choice and control over their journeys, whether on the go or on digital channels.

Read more

Cathay Pacific

Customers are shifting from "do it with me" to more proactive "do it for me" experiences, and AI is a critical enabler in how we provide those experiences, whenever and wherever they need us. As a result, Walmart has been developing conversational AI capabilities for the past several years.

Over time, it has evolved into a robust multi-layered solution that spans across several areas of business to support our customers, members, associates, drivers and marketplace sellers, among others.

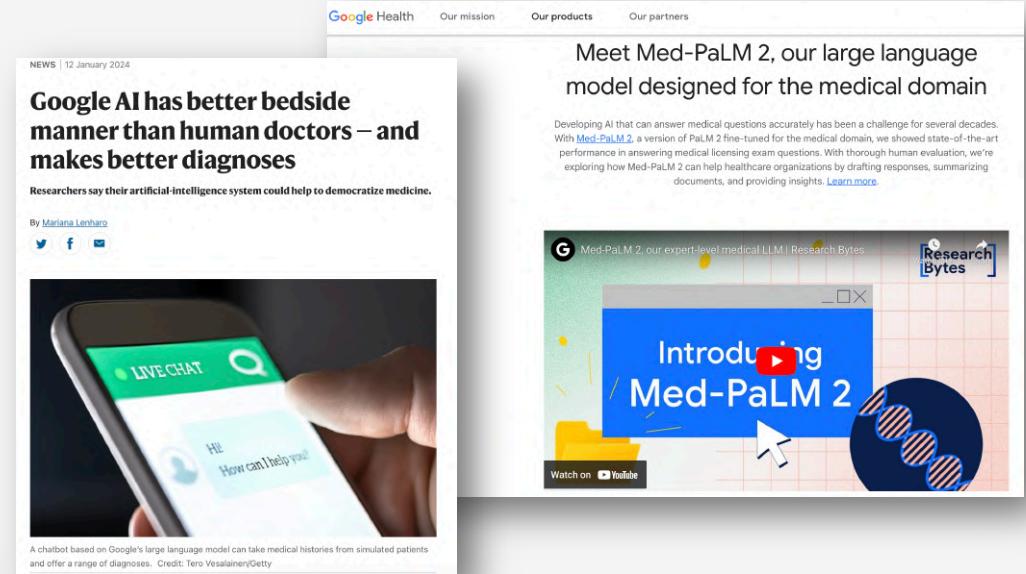
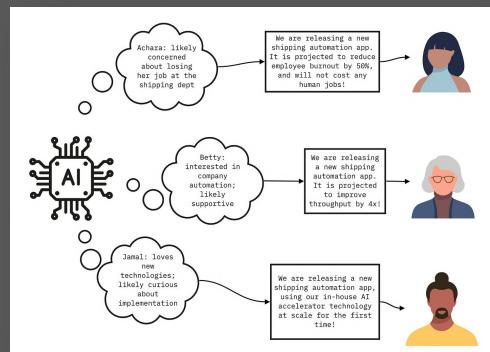
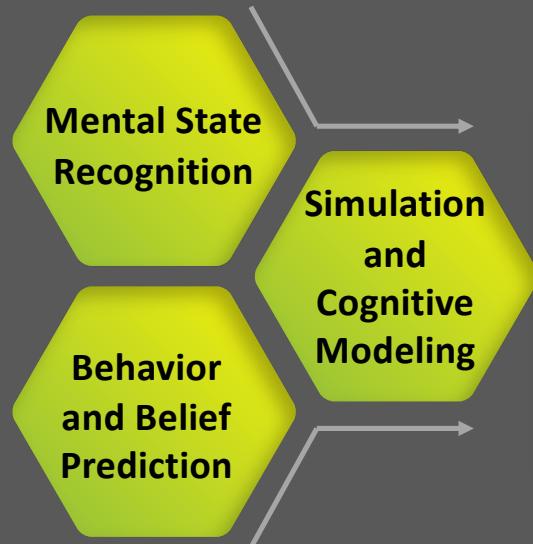
Here's a look at how we're leveraging conversational AI to help our customers and associates around the world save time and have a better experience:

Theory of Mind

01

Introduction

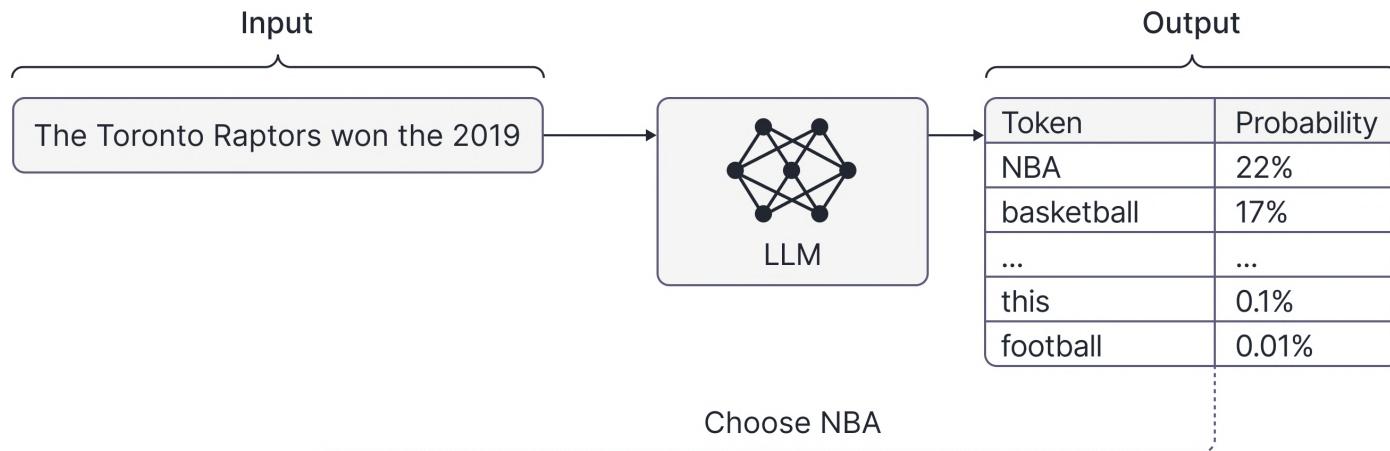
Premack and Woodruff (1978) coined the concept and define **Theory of Mind (ToM)** as the ability that "... one infers states that are not directly observable and one uses these states anticipatorily, to predict the behavior of others as well as one's own." (p. 525)



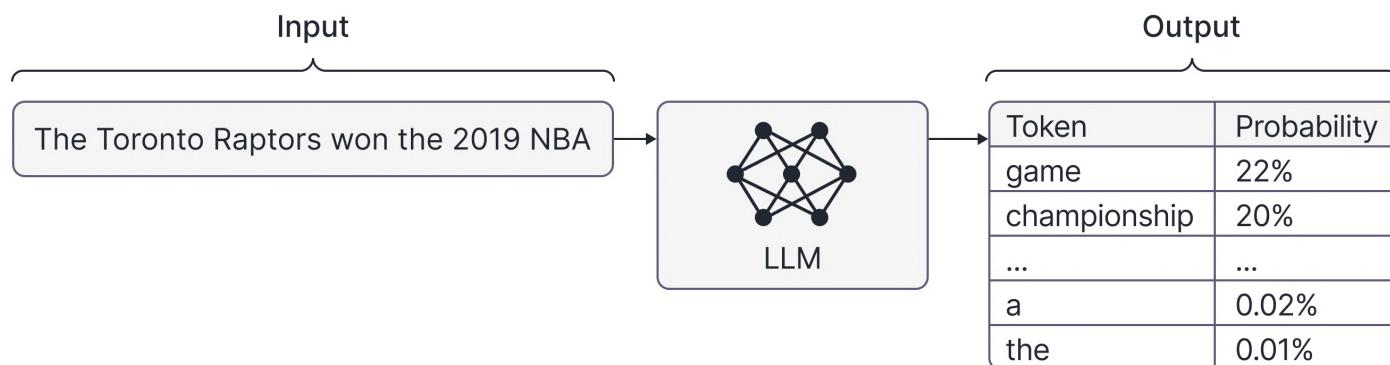
- Has current conversational AI (LLM-empowered) achieved ToM?
- How to equip conversational AI (LLM-empowered) with ToM?

Overview of a General LLM Processing Pipeline

(a) Step N

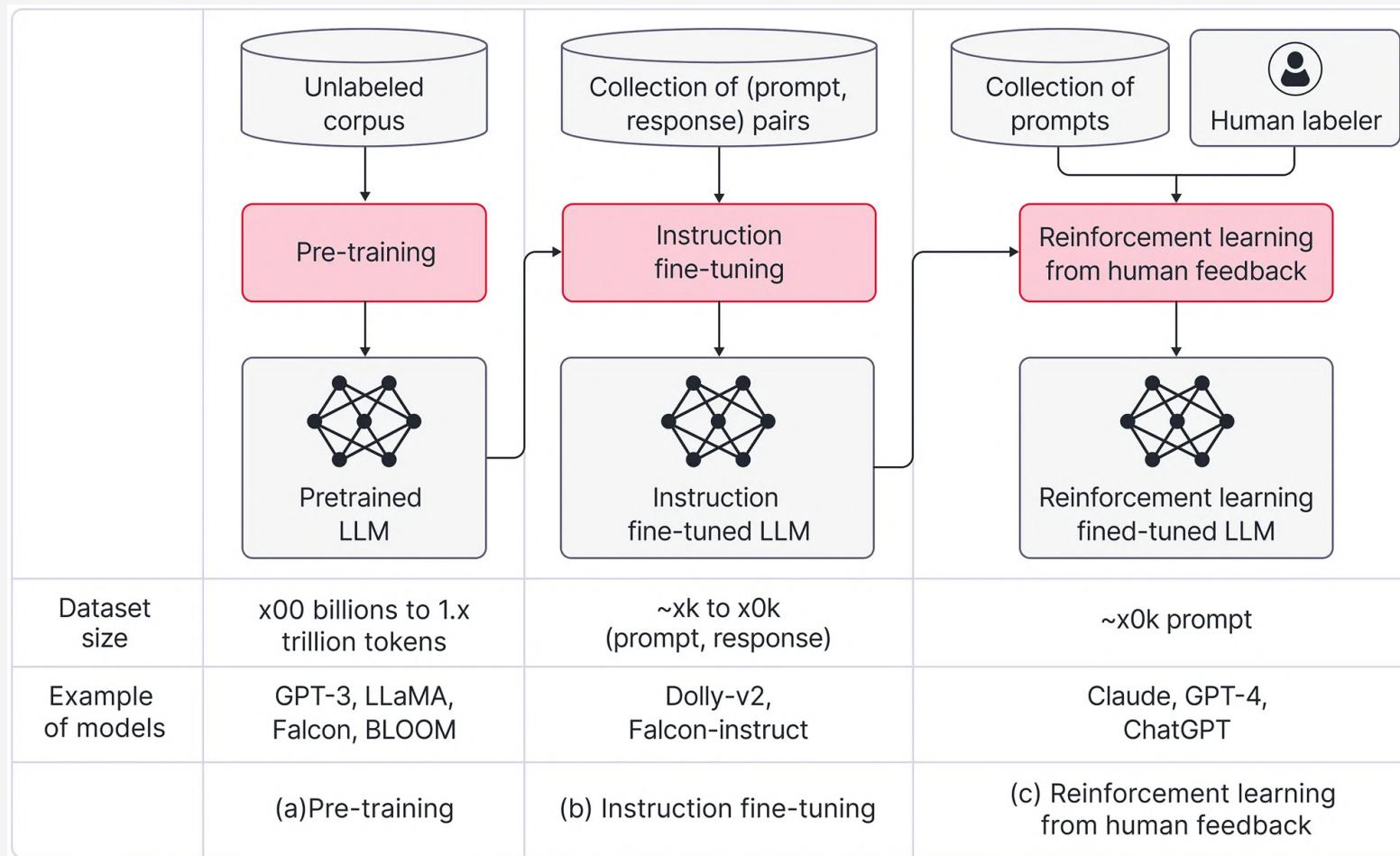


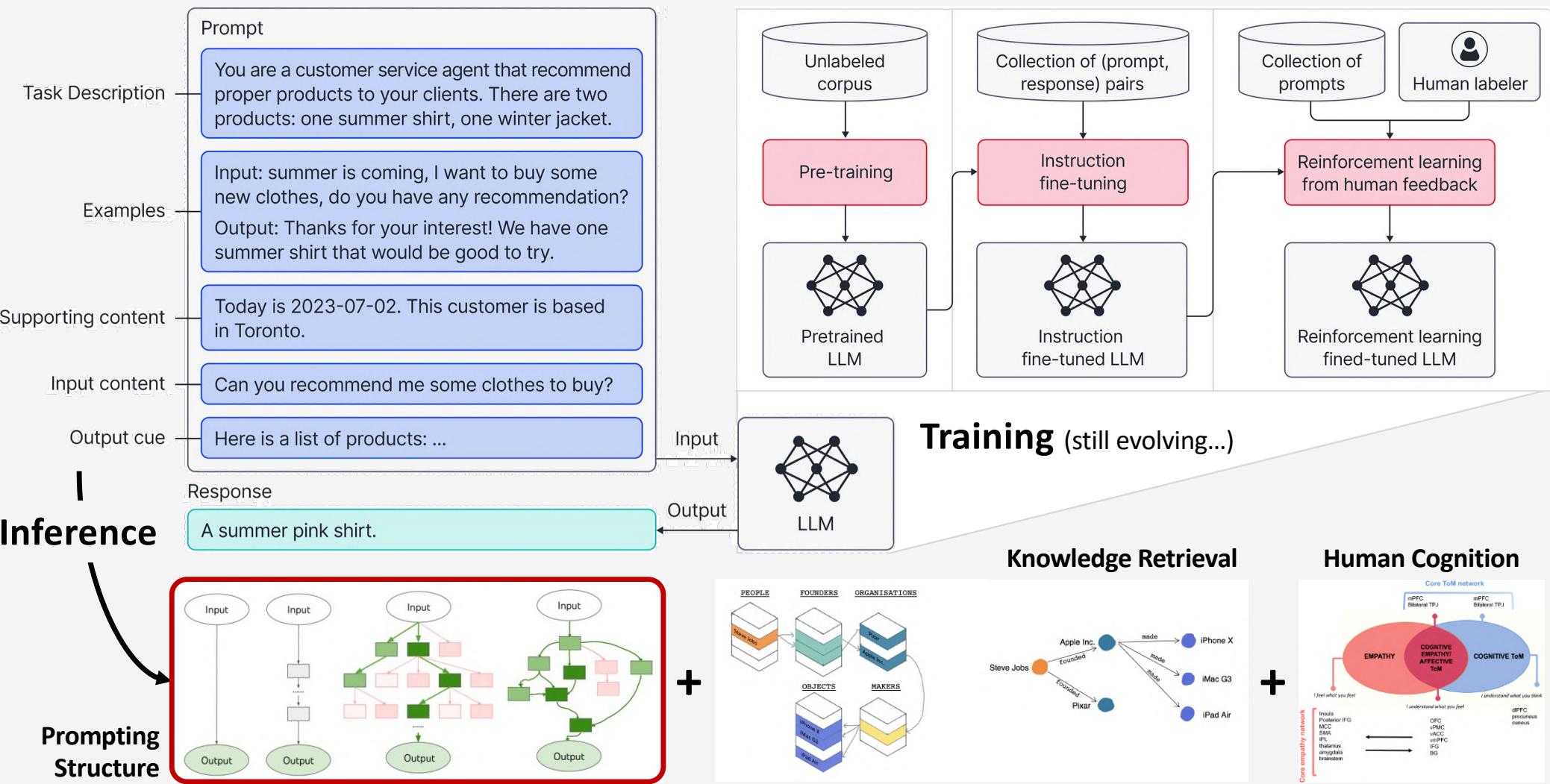
(b) Step N+1



<https://www.borealisai.com/research-blogs/a-high-level-overview-of-large-language-models/>

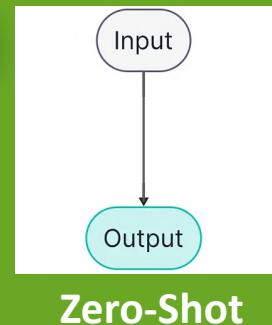
Overview of a General LLM Processing Pipeline





<https://www.borealisai.com/research-blogs/a-high-level-overview-of-large-language-models/>
<https://thesciencemuseum.github.io/heritageconnector/post/2020/11/06/knowledge-graphs-machine-learning-and-heritage-collections/>
https://www.linkedin.com/posts/tonyseale_gpt4-promptengineering-semanticweb-activity-7075381524631580672-TAv3
https://www.researchgate.net/publication/356390737_Social_cognition_in_the_FTLD_spectrum_evidence_from_MRI/figures?lo=1&utm_source=google&utm_medium=organic

Prompt Engineering Techniques (Prompting Paradigms)



Standard Prompting



Model Input

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

Supporting Content +
Input Content

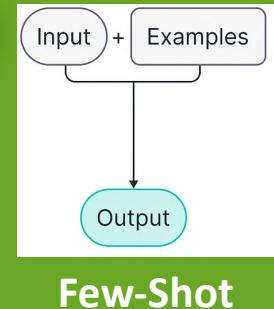
Model Output

A: The answer is 27. ❌

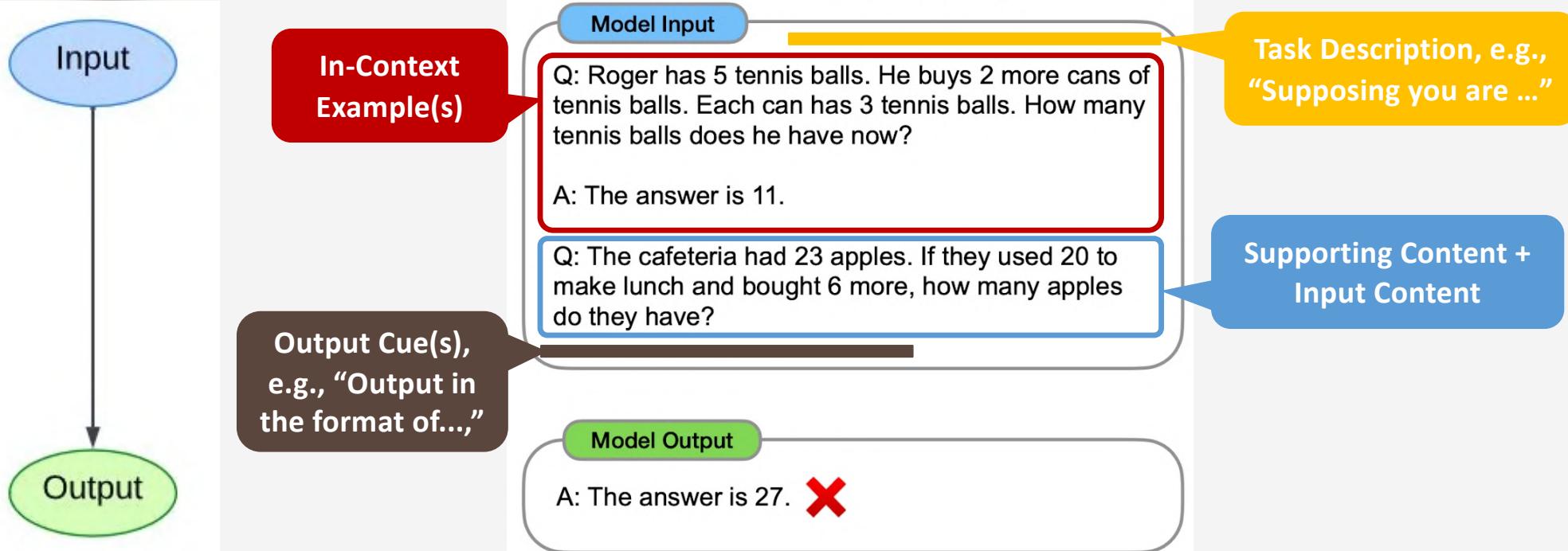


https://www.linkedin.com/posts/tonyseale_gpt4-promptengineering-semanticweb-activity-7075381524631580672-TAv3
<https://deepgram.com/learn/chain-of-thought-prompting-guide>
<https://arxiv.org/abs/2201.11903>

Prompt Engineering Techniques (Prompting Paradigms)

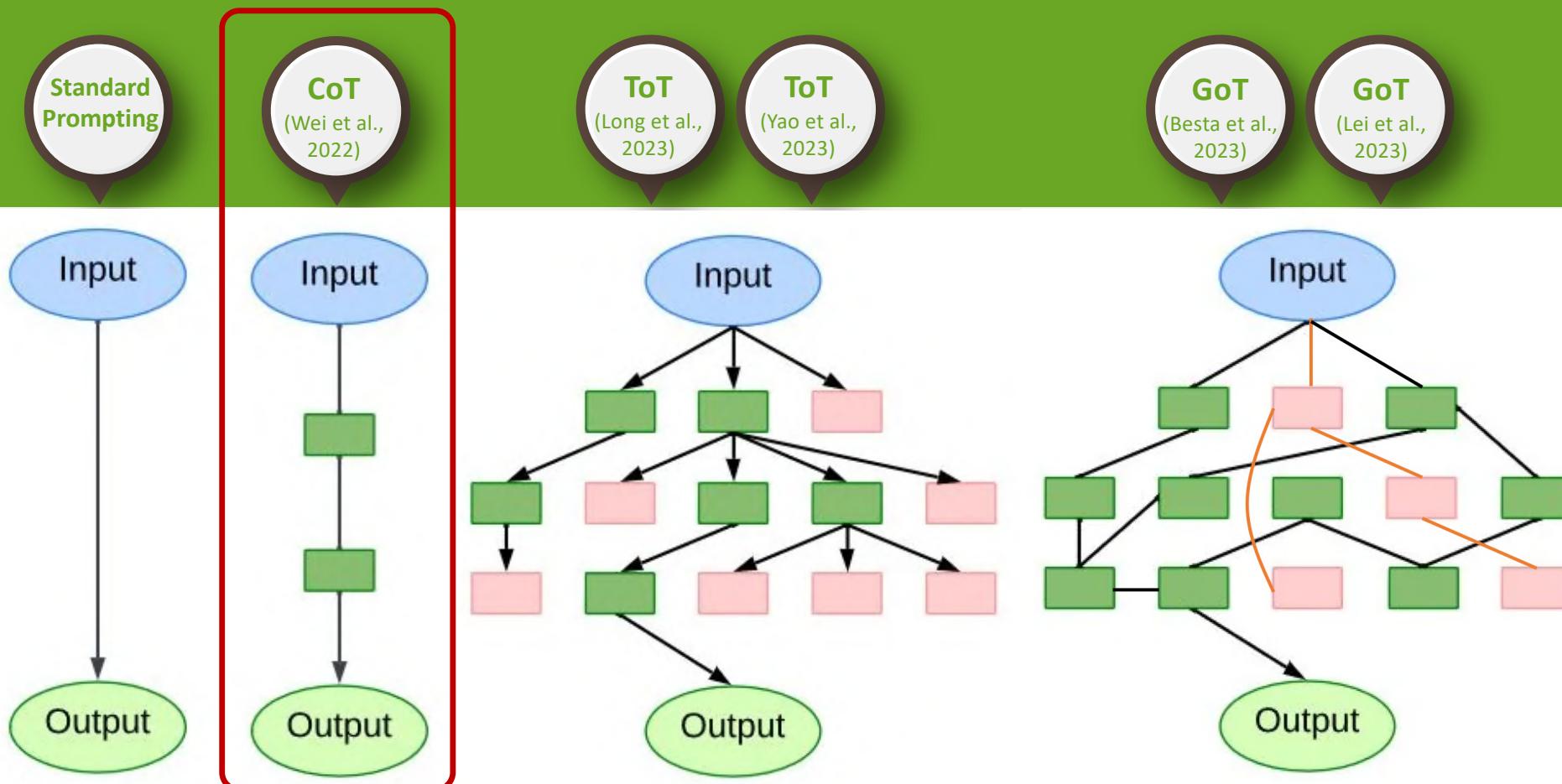


Standard Prompting



https://www.linkedin.com/posts/tonyseale_gpt4-promptengineering-semanticweb-activity-7075381524631580672-TAv3
<https://deepgram.com/learn/chain-of-thought-prompting-guide>
<https://arxiv.org/abs/2201.11903>

Prompt Engineering Techniques (Prompting Paradigms)



O2

Transcending Input-Output Prompts

Chain-of-Thought (CoT)

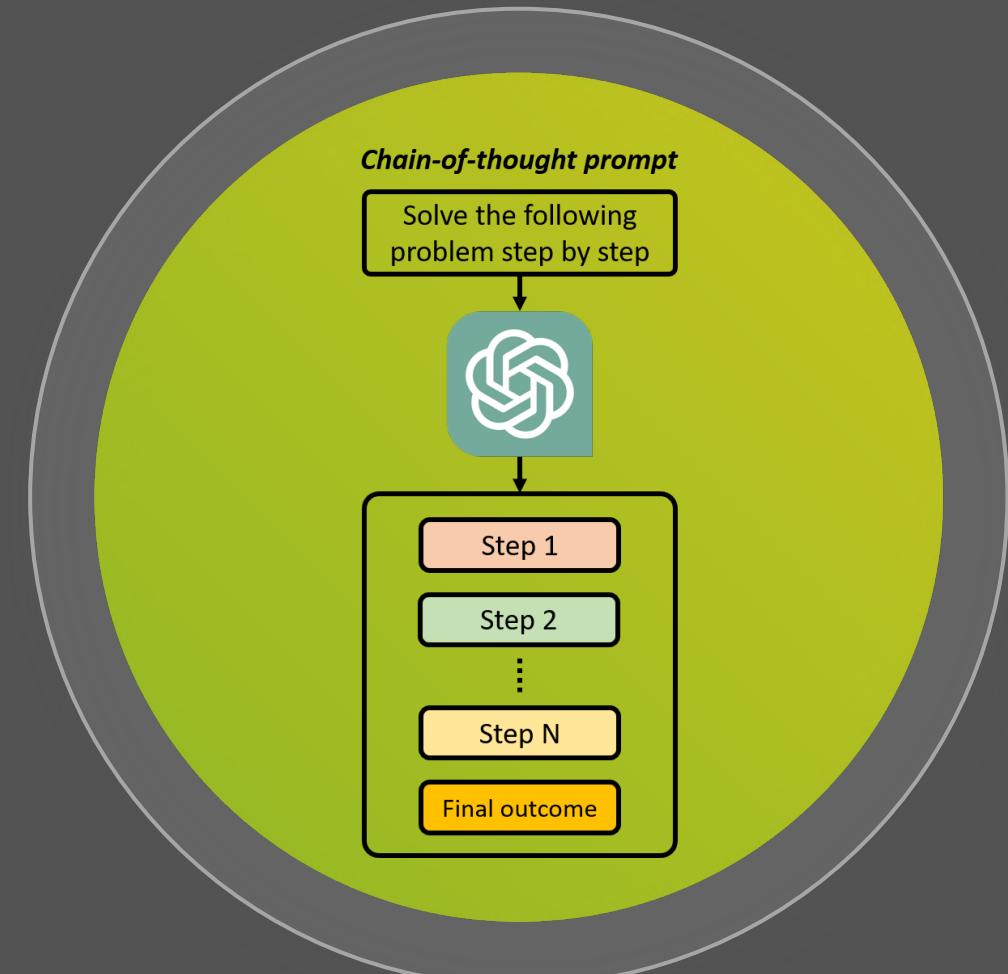
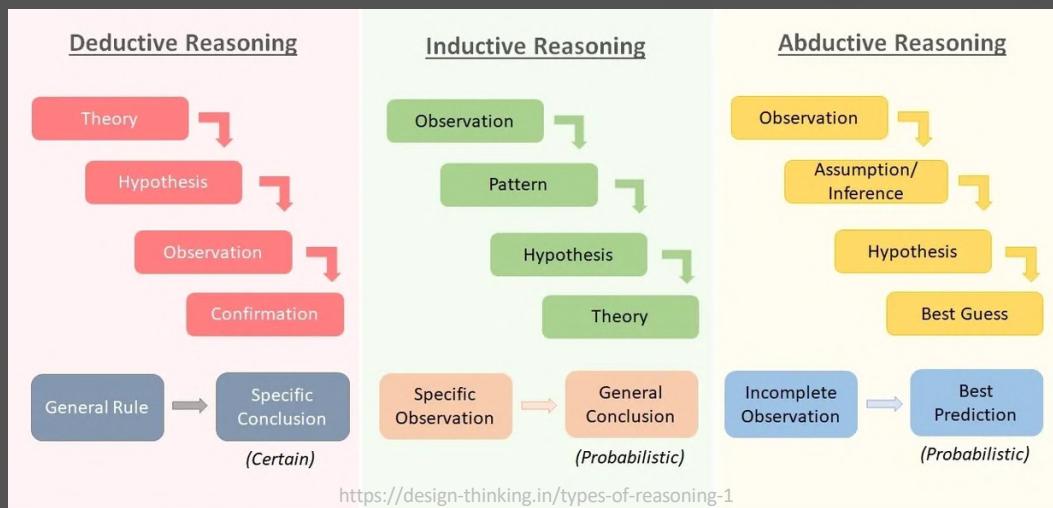
- Introduction of how a chain-of-thought works
- Variations of simple and complex chains-of-thought
- Pros and cons of chains-of-thought prompting

“A chain of thought is a series of intermediate natural language reasoning steps that lead to the final output...” (Wei et al., 2022)

<Input, CoT, Output>

O2

Transcending Input-Output Prompts Chain-of-Thought (CoT)



Components of a CoT Rationale

- **Bridging Objects:** the symbolic items which the model traverses to reach a final conclusion. Bridging can be made up of numbers and equations for arithmetic tasks, or the names of entities in factual tasks.

Arithmetic Reasoning	Multi-hop QA
<p>Q: Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?</p> <p>A: Originally, Leah had 32 chocolates and her sister had 42. So in total they had $32 + 42 = 74$. After eating 35, they had $74 - 35 = 39$ pieces left in total. The answer is 39.</p>	<p>Q: Who is the grandchild of Dambar Shah?</p> <p>A: Dambar Shah (? - 1645) was the father of Krishna Shah. Rudra Shah was the child of Krishna Shah (? - 1661). So the final answer (the name of the grandchild) is: Rudra Shah.</p>

- **Language Templates:** the textual hints (relations / predicates) that guide the language model to derive and contextualize the correct bridging objects during the reasoning process.

“...being relevant [**relevance**] to the query and correctly ordering [**coherence**] the reasoning steps are the key for the effectiveness of CoT prompting” (Wang et al., 2023)

(Wang et al., 2023) <https://aclanthology.org/2023.acl-long.153.pdf>



"Originally, Leah had 32 chocolates and her sister had 42. So in total they had $32 + 42 = 74$ chocolates. So in total they had 74 - 35 = 39."

Coherence

"Q: Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?"

"A: Originally, Leah had 32 chocolates and her sister had 42.

So in total they had $32 + 42 = 74$.

Bridging Objects

After eating 35, they had $74 - 35 = 39$ pieces left in total.

The answer is 39."

Language Templates

"After eating 32, they had 42 pieces left in total. Originally, Leah had 32 chocolates and her sister had 42."

Relevance

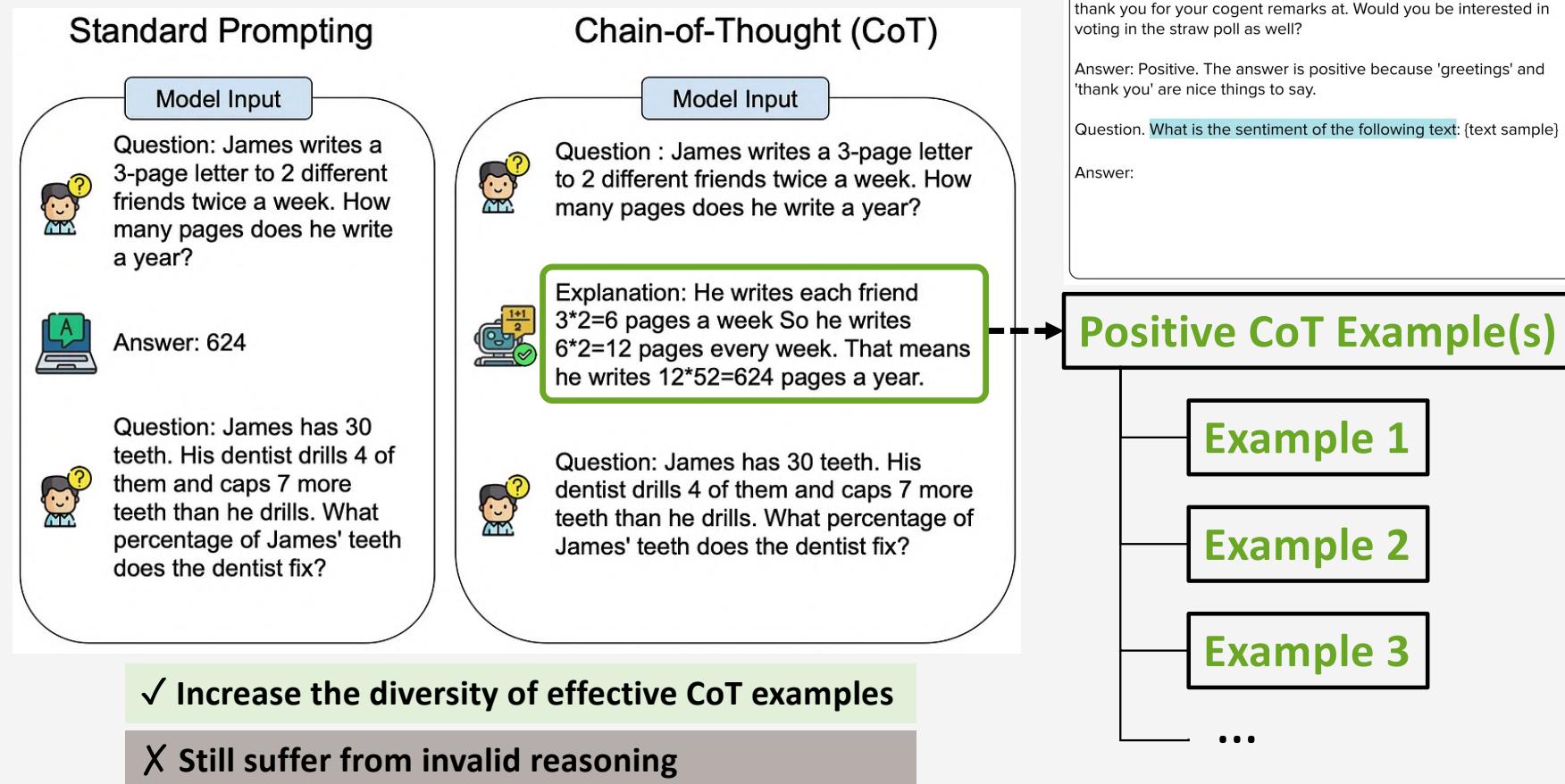
"Originally, Leah had 19 chocolates and her sister had 31. So in total they had $19 + 31 = 50$."

"Patricia needs to donate 32 inches and wants her hair to be 42 inches long after donating."

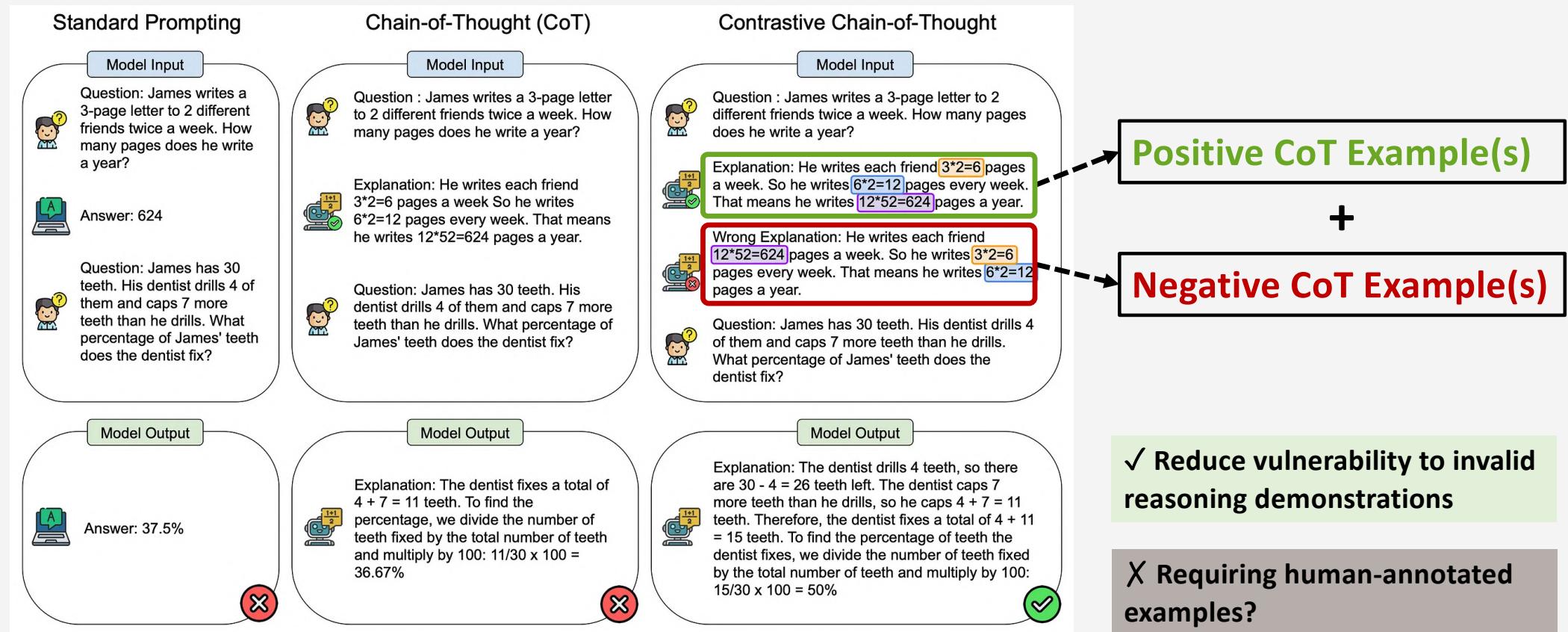
Relevance matters more than coherence for bridging objects and Coherence of language templates is important.



Few-Shot Chain-of-Thought Prompting



Contrastive Chain-of-Thought Prompting



Zero-Shot Chain-of-Thought Prompting

“Let’s think step by step.”



(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. ✗

(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 ✗

(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. ✓

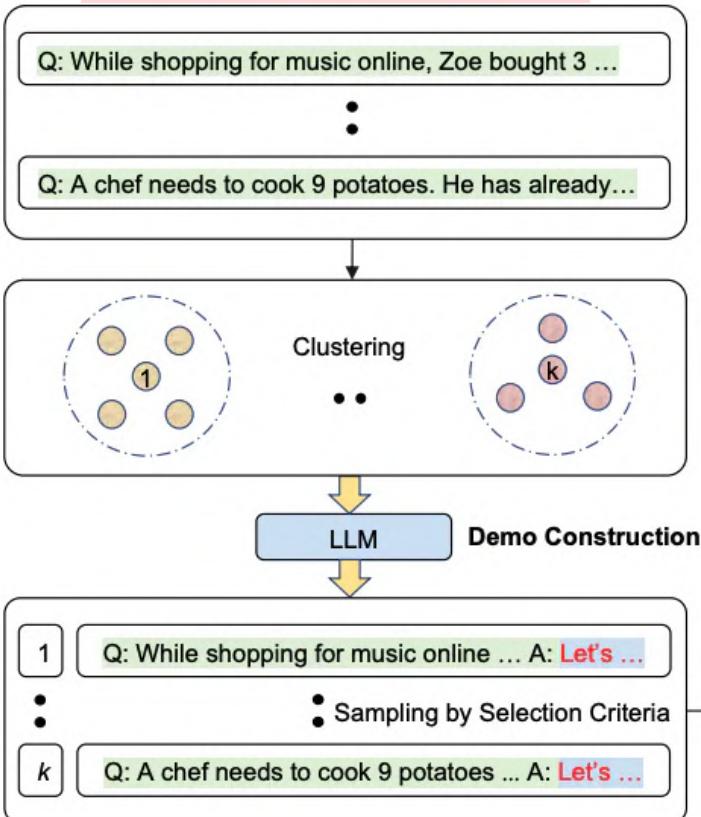
✓ Mitigate
the trouble
of coming
up with
(too many)
examples

✗ May not
work well
in all cases



Automatic Chain-of-Thought Prompting (Auto-CoT)

Stage 1): question clustering



Stage 2): demonstration sampling

Auto Demos One by One

Zero-shot CoT

Q: While shopping for music online, Zoe bought 3 country albums and 5 pop albums. Each album came with a lyric sheet and had 3 songs. How many songs did Zoe buy total?

A: **Let's think step by step.** Zoe bought 3 country albums. Each album has 3 songs. So she bought $3 \times 3 = 9$ songs from the country albums. Zoe bought 5 pop albums. Each album has 3 songs. So she bought $5 \times 3 = 15$ songs from the pop albums. Zoe bought $9 + 15 = 24$ songs in total. The answer is 24.
...

Q: A chef needs to cook 9 potatoes. He has already cooked 7. If each potato takes 3 minutes to cook, how long will it take him to cook the rest?

A: **Let's think step by step.** The chef has already cooked 7 potatoes. That means it has taken him 7×3 minutes to cook those 7 potatoes. That means it will take him 3 more minutes to cook each of the remaining 2 potatoes ...

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?

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

Test Question

LLM In-Context Reasoning

The pet store had 64 puppies. They sold 28 of them. That means they have 36 puppies left. They put the rest into cages with 4 in each cage. That means they have 9 cages. The answer is 9.

Inference

✓ Auto-generation of positive (and negative) examples

✗ Depend on the outputs of Stage 1&2

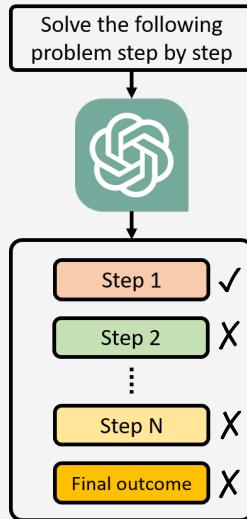


Self-Consistency Chain-of-Thought Prompting

Manually or through
special prompt, e.g.,

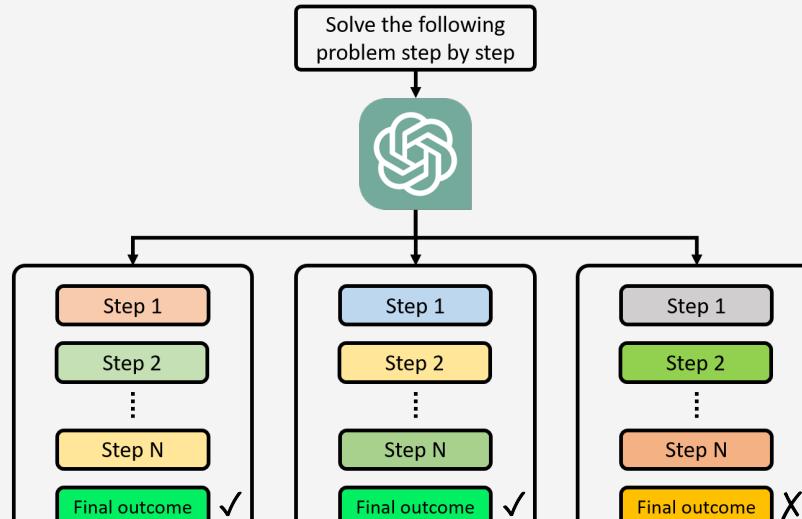
"Imaging X independent experts who reason differently..."

Chain-of-thought prompt

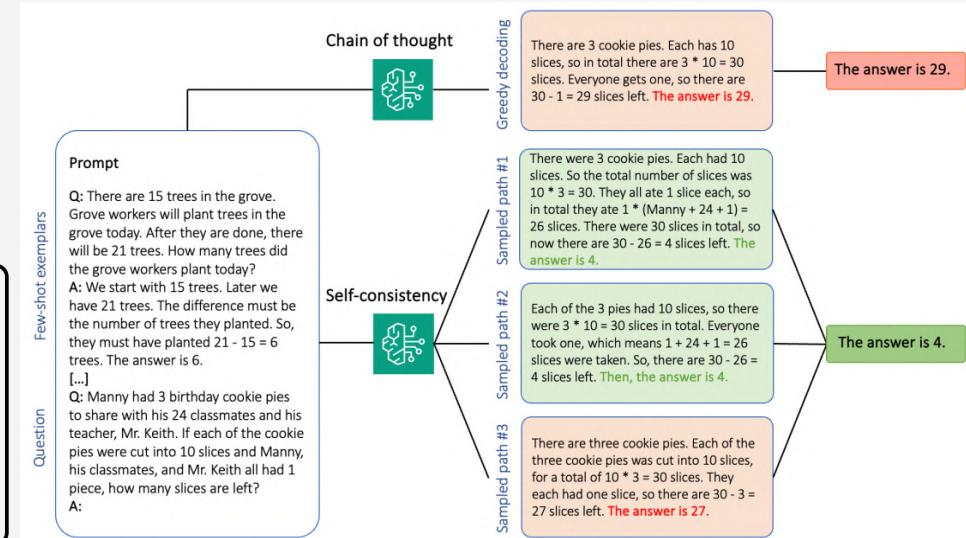
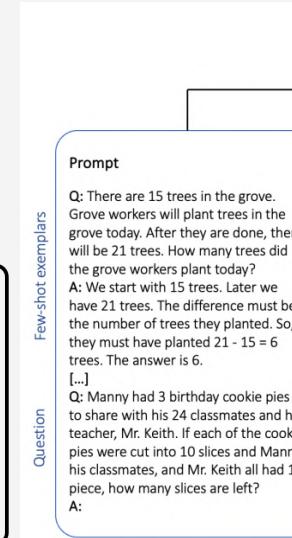


A single CoT

Self-consistency prompt



Majority Vote



✓ Improve the consistency of an existing LLM without altering its architecture

✓ May leverage different LLMs and multi-agent setup

✗ Require large output tokens



(Wang et al. 2022) <https://arxiv.org/abs/2203.11171>

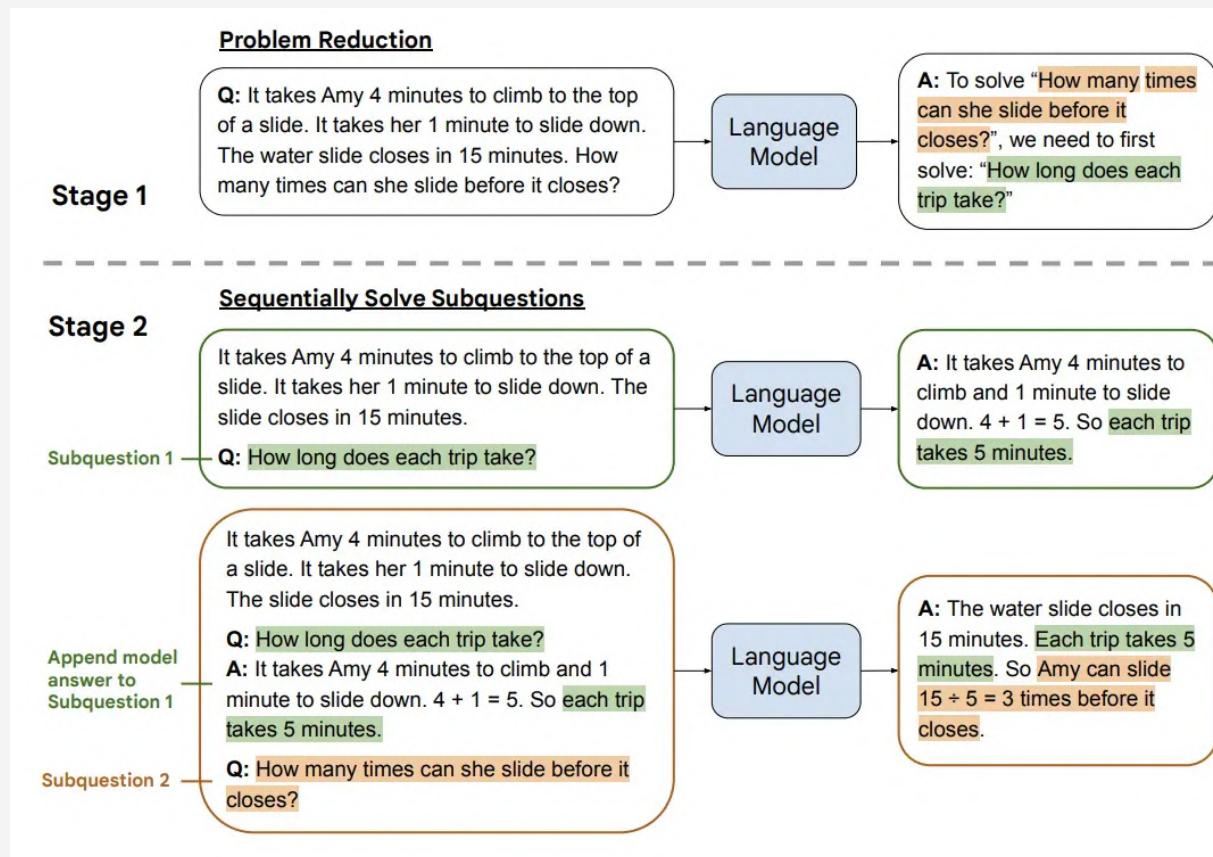
<https://campus.datacamp.com/courses/chatgpt-prompt-engineering-for-developers/advanced-prompt-engineering-strategies?ex=10>

<https://aws.amazon.com/blogs/machine-learning/enhance-performance-of-generative-language-models-with-self-consistency-prompts-on-amazon-bedrock/>

Least to Most (Chain-of-Thought) Prompting

"What subproblems must be solved before answering the inquiry?."

These subproblems are solved one by one. The solution of previous subproblems is fed into the prompt trying to solve the next problem.

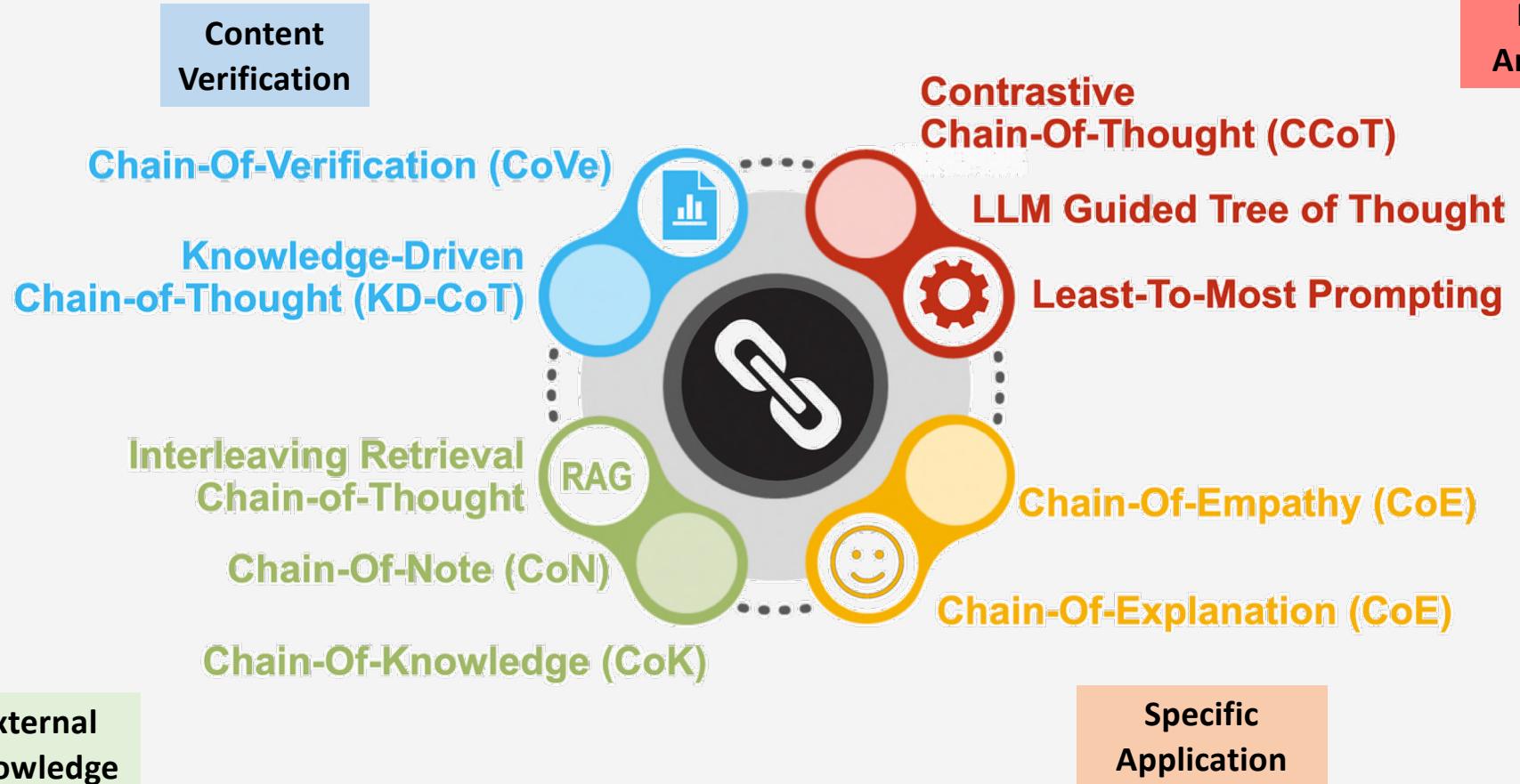


✓ Mitigate the trouble of coming up with (too many) examples

X May not work well in all cases



Existing (and Ever-Growing) Family of Chain-of-X Prompting



CoT Vulnerabilities

- Hallucination
- Error propagation
- Lack of external knowledge
- “CoT only yields performance gains when used with models of $\sim 100B$ parameters”

CoT Key Aspects

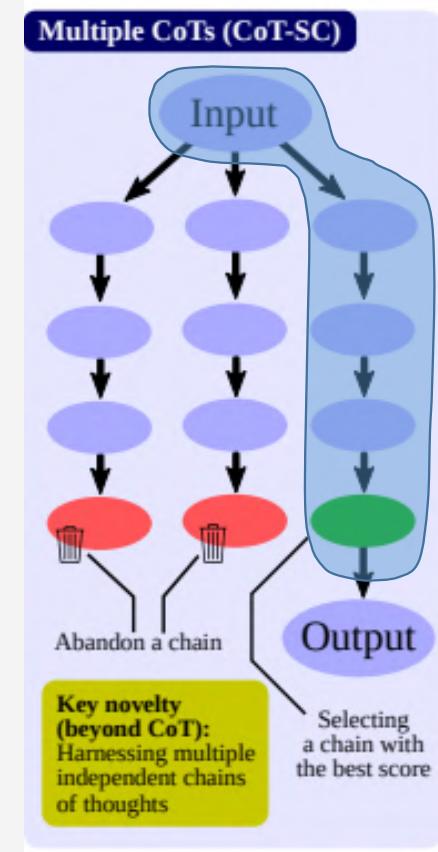
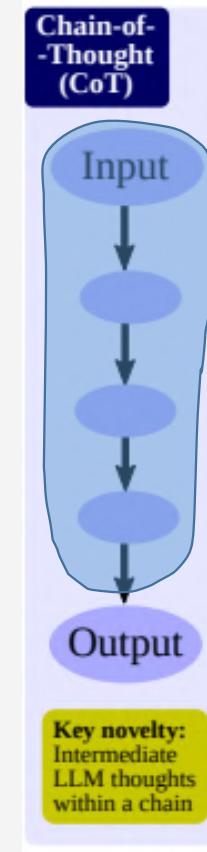
- Self-consistency
- Sensitivity to prompt design
- Coherence of the order of steps in a CoT rationale

Assuming a **fixed thought size** (#tokens) and a **fixed context size N** (#thoughts in the LLM context)

Latency: Number of steps between the input and output



Volume: For a thought t , the number of preceding thoughts that could impact t

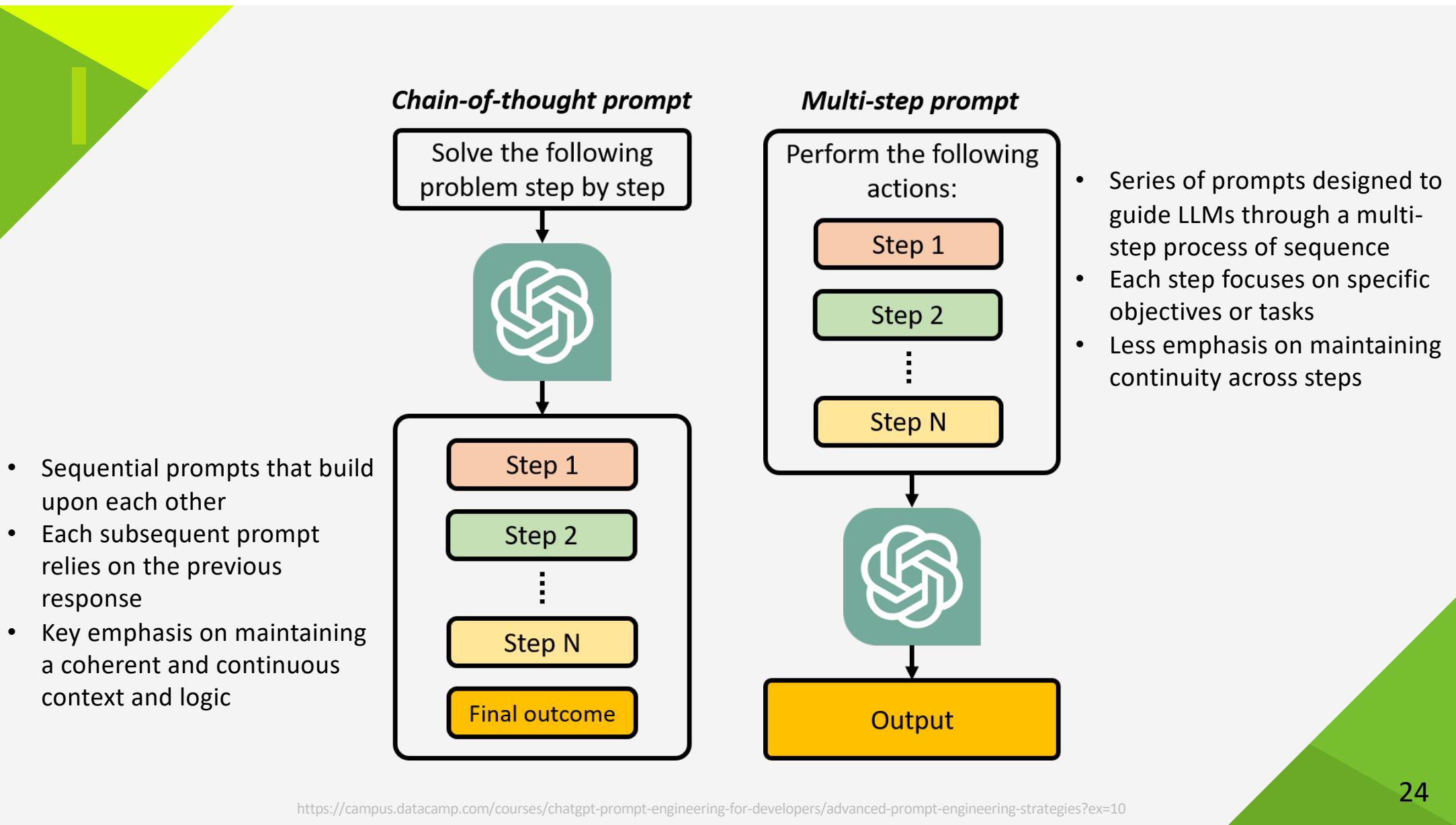


✓✓ Large Volume: N

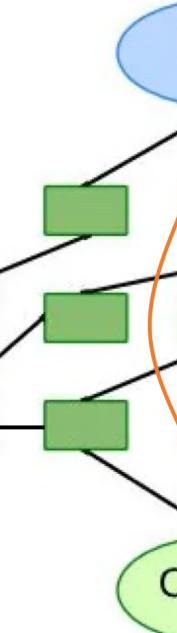
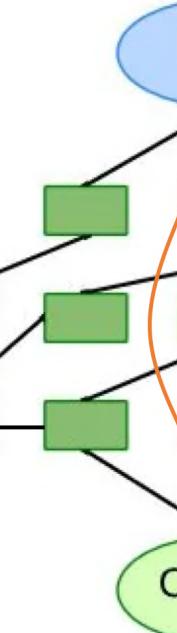
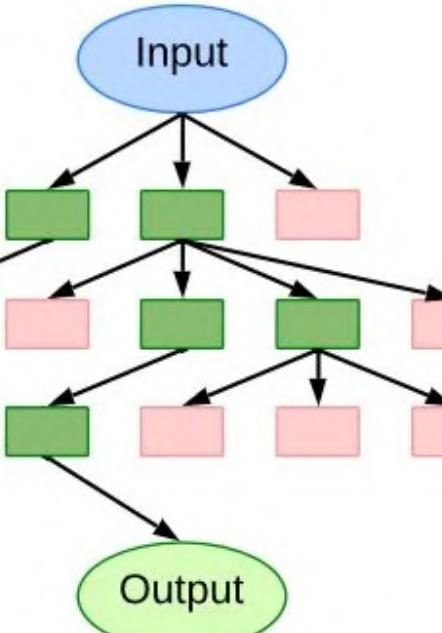
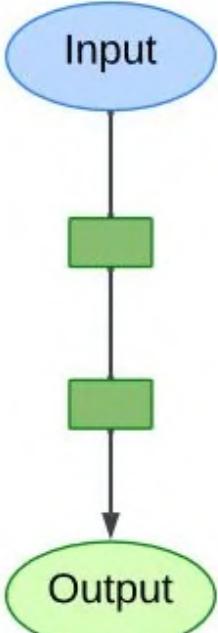
✗✗ Large Latency: N

✗ Low Volume: N/k

✓ Low Latency: N/k



Prompt Engineering Techniques (Prompting Paradigms)



“LMs ... might benefit from augmentation by a more deliberate ‘System 2’ planning process that ... explores diverse alternatives ... and evaluates its current state” (Yao et al., 2023)

03

Deliberate Problem Solving Tree-of-Thoughts

- Description of how a tree-of-thoughts works
- Concepts and mechanisms of trees-of-thoughts
- Pros and cons of using trees-of-thoughts

System 1
<ul style="list-style-type: none">• Fast• Automatic• Non-conscious• Common Use• ‘Good Enough’

System 2
<ul style="list-style-type: none">• Slow• Deliberative• Conscious• Requires Effort• Use If Critical

03

Deliberate Problem Solving Tree-of-Thoughts

Decomposition

Breaks a problem into a sequence of smaller steps—or thoughts—that are solved individually

Generation

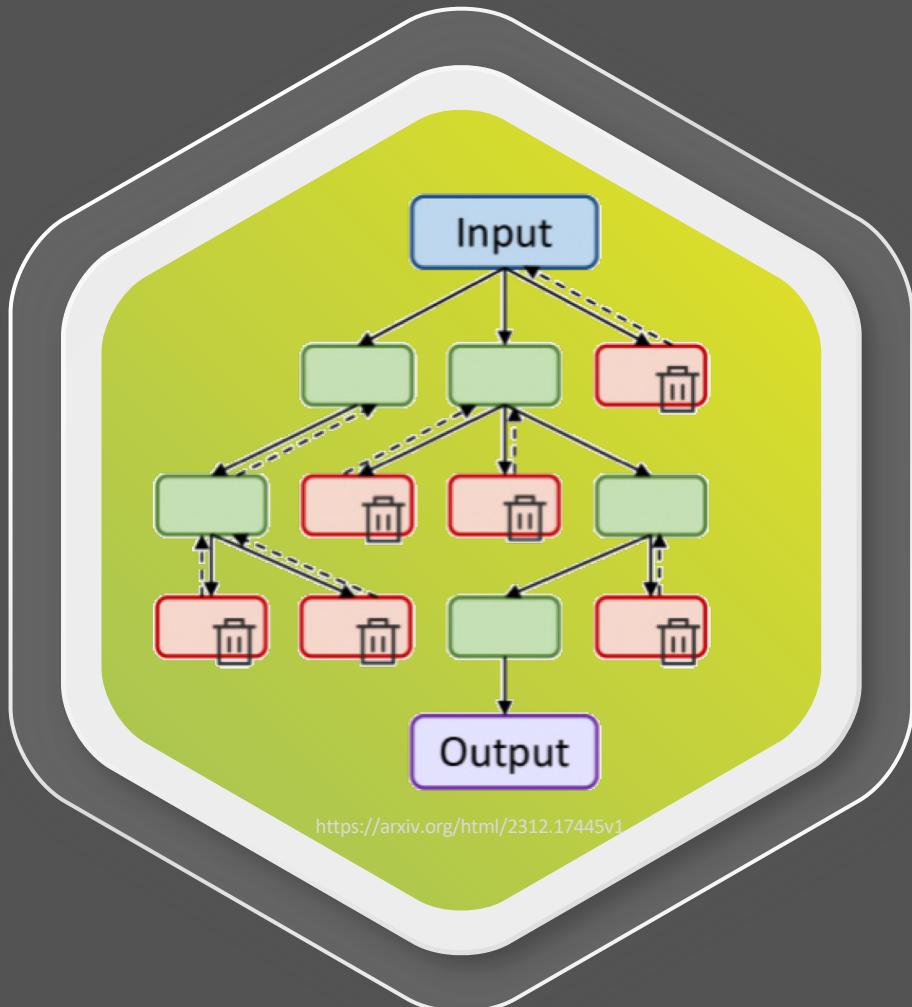
Explore multiple choices for each problem-solving thought.

Evaluation

Evaluate whether certain thoughts bring the model closer to a final solution

Search

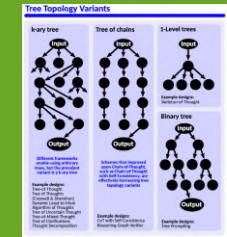
Search over a combinatorial space of possible problem-solving steps to find the best final solution. Perform backtracking when certain thoughts are found to be a dead end.



1)

Thought Decomposition: Break Intermediate Process into Steps

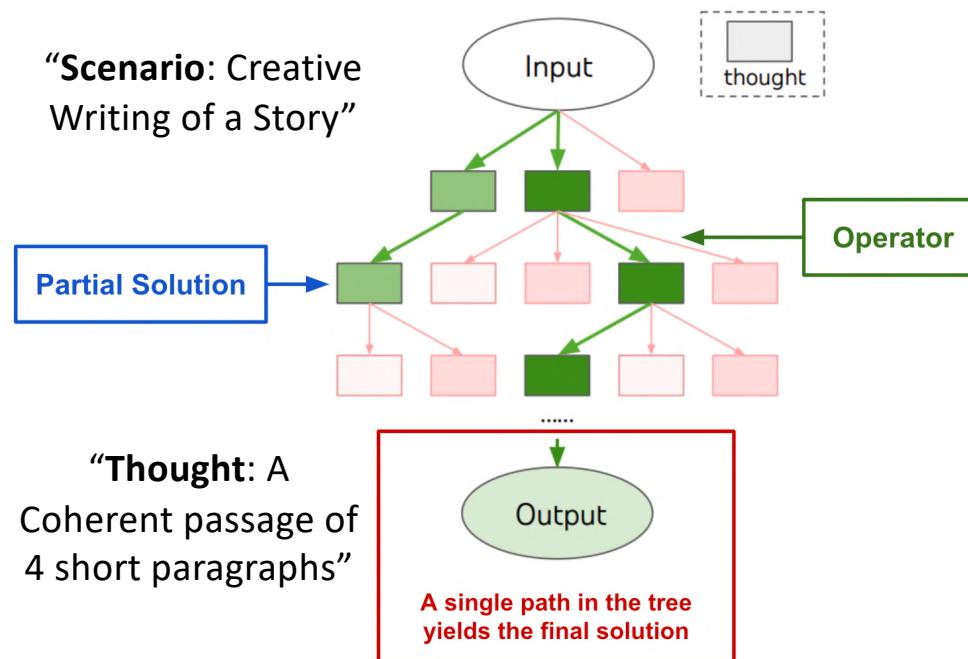
- Definition of “thought” is task dependent and may take different forms



“A thought ... should be **small enough** so that LMs can generate **promising and diverse samples** ...”

E.g., generating a whole book is usually too big to be coherent

“Scenario: Creative Writing of a Story”



“A thought ... should be ... yet **big enough** so that LMs can **evaluate its prospect** toward problem solving ...”

E.g., generating one word in a paragraph is usually too small to evaluate

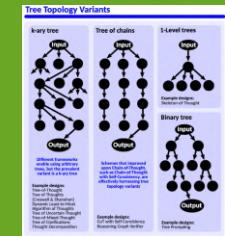


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

1)

Thought Decomposition: Break Intermediate Process into Steps

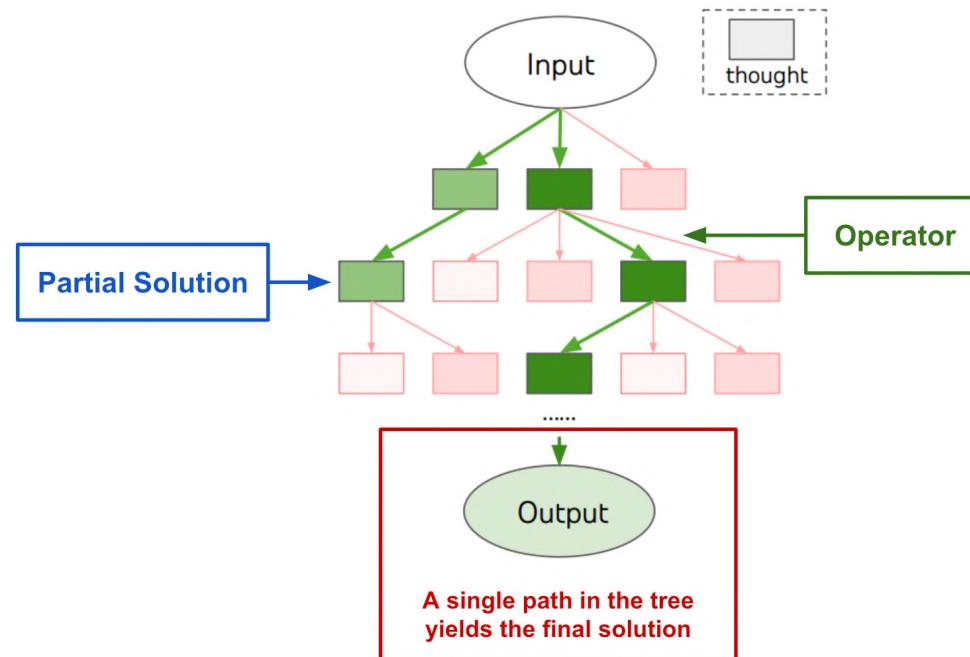
- Definition of “thought” is task dependent and may take different forms
- The “thoughts” in the tree can be homogeneous or heterogeneous
- The topology of the tree (depth, degrees, etc.) may vary based on the decomposition



“Scenario: Sudoku”

- Insert a number
- Insert a number
- Insert a number

Each level has the same purpose; the number of sub-branches is determined by the rule.



“Scenario: Design a marketing strategy for a new product launch”

- Identify target audience
- Evaluate competitor and trend
- Design promotional campaign

Each level has a different purpose; the number of sub-branches under each node is not fixed and can be pre-defined.

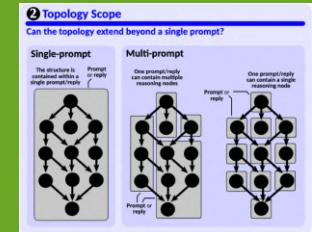


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f>

2)

Thought Generation: Generate K Candidates Given a Tree State

- Definition of “thought” is task dependent and may take different forms
- **Sampling:** generating several thoughts independently with the same prompt
- **Proposing:** generating several thoughts sequentially with a “propose prompt”

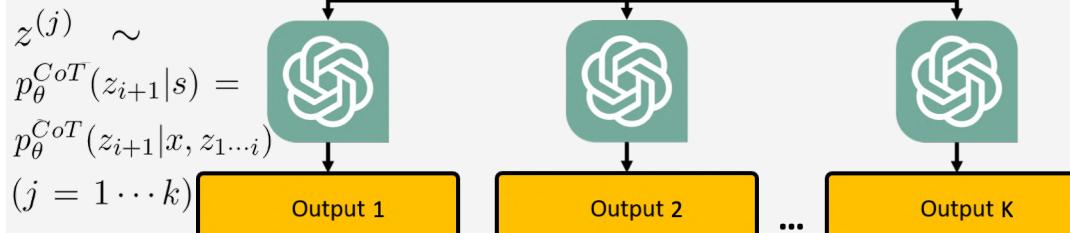


The sampling approach works best when the thought space is rich, as several independently-generated thoughts are unlikely to be duplicates.

Thought Sampling

Perform the following actions:
Step 1
Step 2
⋮
Step N

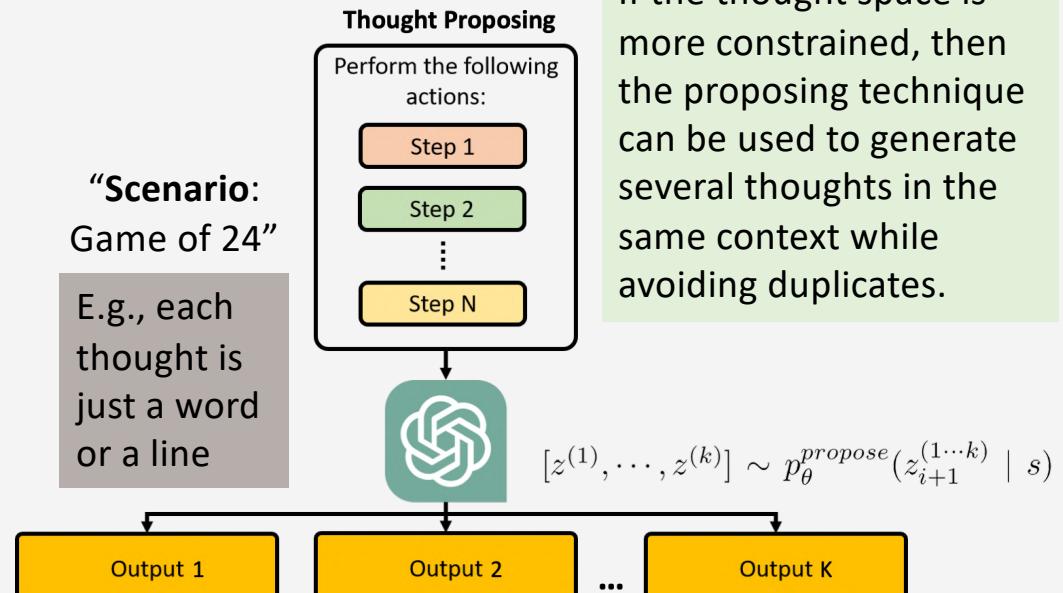
“Scenario:
Creative Writing
E.g., each thought is a paragraph



Thought Proposing

Perform the following actions:
Step 1
Step 2
⋮
Step N

“Scenario:
Game of 24
E.g., each thought is just a word or a line



If the thought space is more constrained, then the proposing technique can be used to generate several thoughts in the same context while avoiding duplicates.

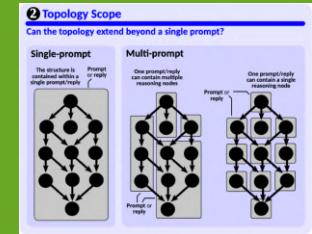


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

2)

Thought Generation: Generate K Candidates Given a Tree State

- Definition of “thought” is task dependent and may take different forms
- **Sampling:** generating several thoughts independently with the same CoT prompt
- **Proposing:** generating several thoughts sequentially with a “propose prompt”

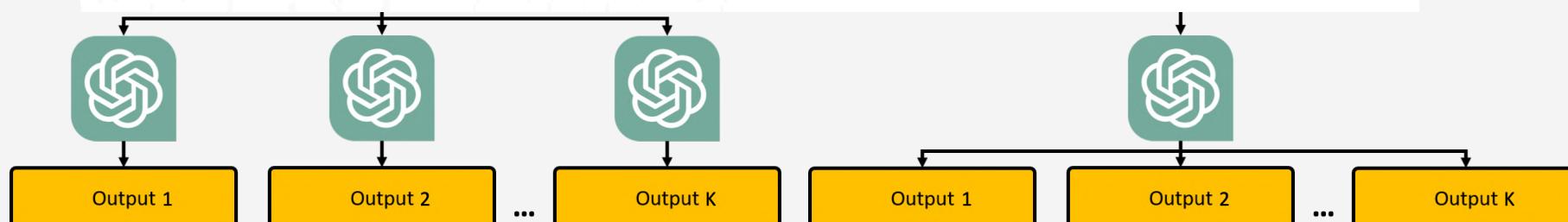


Thought Sampling

Sampling or Proposing?

Thought Proposing

Prompt: *Imagine three different experts are answering this question. All experts will write down 1 step of their thinking, then share it with the group. Then all experts will go on to the next step, etc. If any expert realizes they're wrong at any point, then they leave. The question is...*

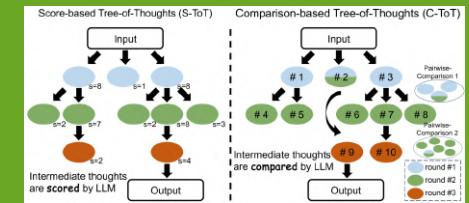


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://www.semanticscholar.org/reader/eff9d7ed0f30f121d30ee13802a11f172ef66f4>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://medium.com/@kyeg/unleashing-super-intelligence-with-tree-of-thoughts-f2f744786e65>

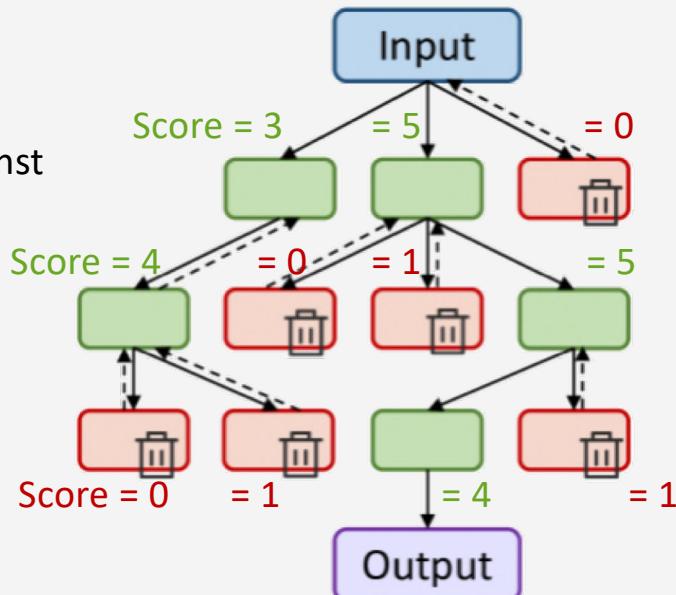
3) State Evaluation: Deliberate Reasoning of Progress to a Solution

- Define a “heuristic” for evaluation: a set of criteria or a holistic measure
- **Value:** independently assign a scalar value or classification to each state



Score-based:

- Value formula
- Violations against given criteria
- Distance to success



1. Write a story featuring Mowgli.
 - *Choice 1: Mowgli is lonely and wants friends.
 - *Choice 2: Mowgli is adventurous and wants to explore the villages.
 - *Choice 3: Mowgli is curious and wants to learn about the human kingdom
2. Pick one choice and keep writing the story.
3. Check your story for how creative, clear, and well-written it is.

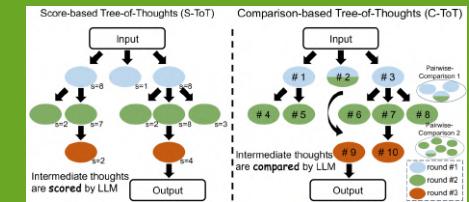


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

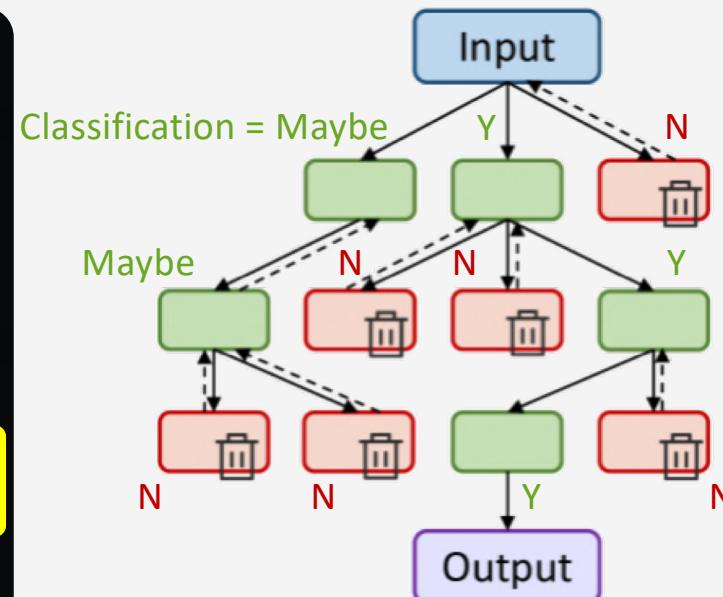
<https://dev.to/zokizuan/tree-of-thoughts-a-new-way-to-prompt-ai-2d1e>

3) State Evaluation: Deliberate Reasoning of Progress to a Solution

- Define a “heuristic” for evaluation: a set of criteria or a holistic measure
- **Value:** independently assign a scalar value or classification to each state
 - Current state only vs. look-ahead simulation
 - Promoting “good” states and/or eliminating “bad” states



1. Write a story featuring Mowgli.
 - *Choice 1: Mowgli is lonely and wants friends.
 - *Choice 2: Mowgli is adventurous and wants to explore the villages.
 - *Choice 3: Mowgli is curious and wants to learn about the human kingdom
2. Pick one choice and keep writing the story.
3. Check your story for how creative, clear, and well-written it is.



Classification-based:

- Categorical classifier
- Ordinal scale that could be heuristically converted into values



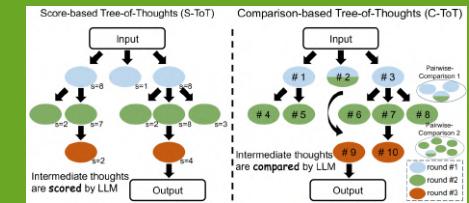
(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://dev.to/zokizuan/tree-of-thoughts-a-new-way-to-prompt-ai-2d1e>

3)

State Evaluation: Deliberate Reasoning of Progress to a Solution

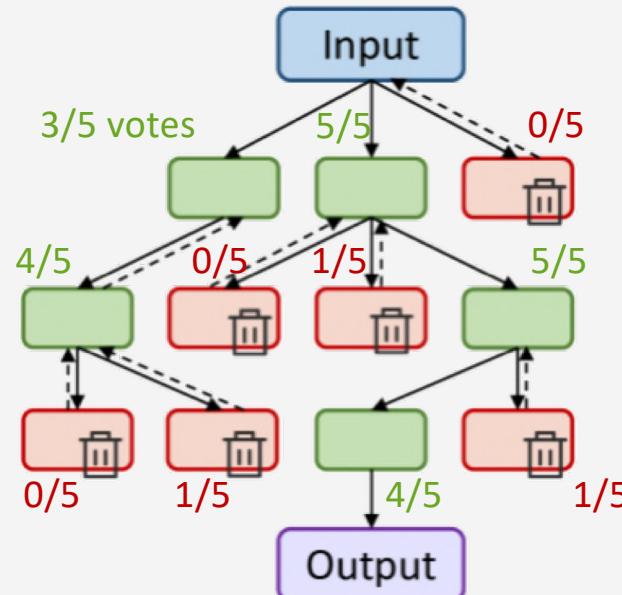
- Define a “heuristic” for evaluation: a set of criteria or a holistic measure
- **Value:** independently assign a scalar value or classification to each state
- **Vote:** compare different solutions and select the one that is most promising



Comparison-based:

- Pair-wise comparison
- Ranking
- Top K
- ...

An example of pair-wise tournament:
 “Compare thoughts a and b by LLM n times based on the given needs, then take a majority vote. If a wins, keep thought a in the pool and drop b , and vice versa.”



1. Imagine you're trying to find a good place to go on vacation.
 Your needs are:
 *Find a good vacation spot. Keep in mind:
 *Must be in India.
 *Good for families.
 *Has fun outdoor activities.
 *Has interesting historical places.
 *Summer weather should be bearable.

2. Choices for looking up places:
 *Choice 1: Look at places in Northern India.
 *More choices: Himachal Pradesh or Uttarakhand.
 *Choice 2: Look at places in Southern India.
 *More choices: Kerala or Tamil Nadu.
 *Choice 3: Look at places in Western India.
 *More choices: Goa or Rajasthan.

3. Pick a choice and dig deeper.
 4. Compare all the places you found based on your needs.

Here, you can use a ‘browser’ tool to help you look for info along the path you chose. The AI can follow a path, do the searching, and check how good the options are based on what you need. This way of solving problems can give you more complete and fitting answers.



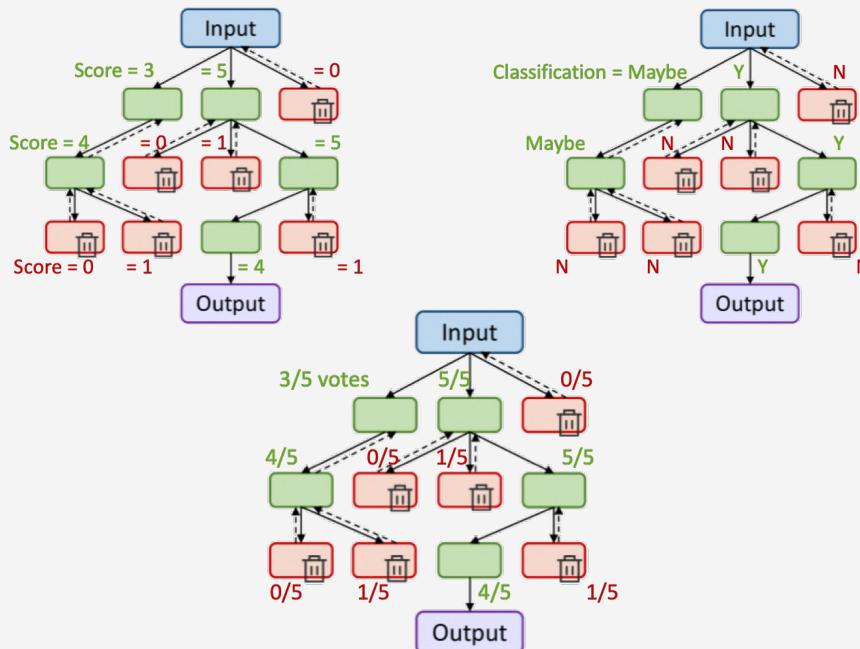
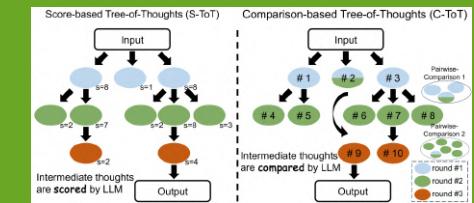
(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://dev.to/zokizuan/tree-of-thoughts-a-new-way-to-prompt-ai-2d1e>

3)

State Evaluation: Deliberate Reasoning of Progress to a Solution

- Define a “heuristic” for evaluation: a set of criteria or a holistic measure
- **Value:** independently assign a scalar value or classification to each state
- **Vote:** compare different solutions and select the one that is most promising



Voting is best when a successful solution to a problem is hard to directly value (e.g., creative writing tasks), e.g., casting “which state to explore” as a multi-choice QA by comparing different partial solutions

For both methods, could prompt the LM multiple times to aggregate the value or vote results to trade time/resource/cost for more faithful/robust heuristics

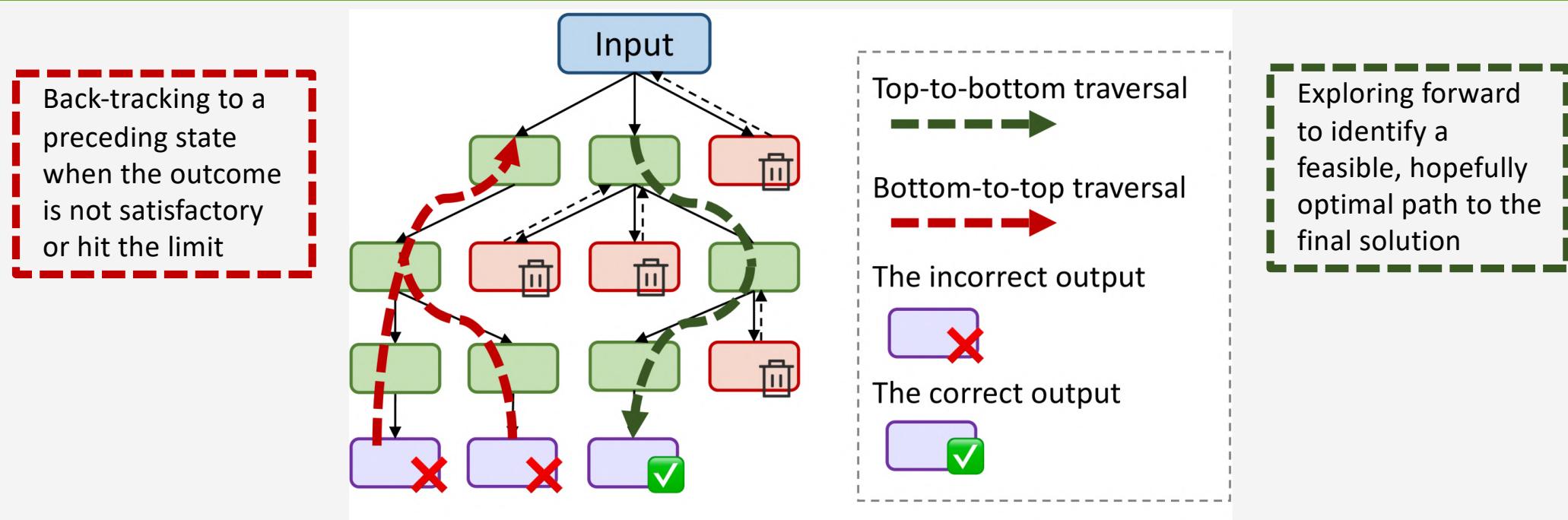
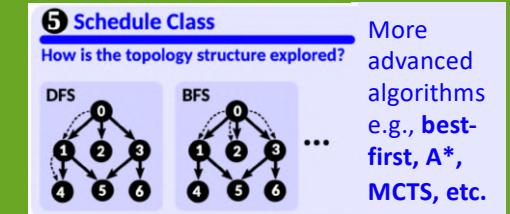


(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

4)

Search Algorithm: Explore the Solution Space

- Plug and play a proper tree search algorithm based on the tree structure



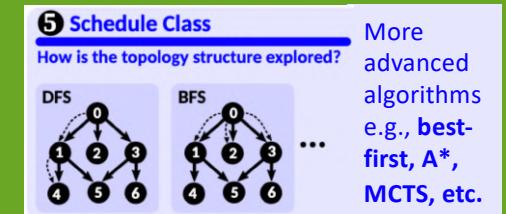
(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>

4)

Search Algorithm: Explore the Solution Space

- Plug and play a proper tree search algorithm based on the tree structure
- **BFS:** maintaining a set of the b most promising states per step
 - Local vs. global evaluation at each layer
 - Set b to a small number to prune the tree



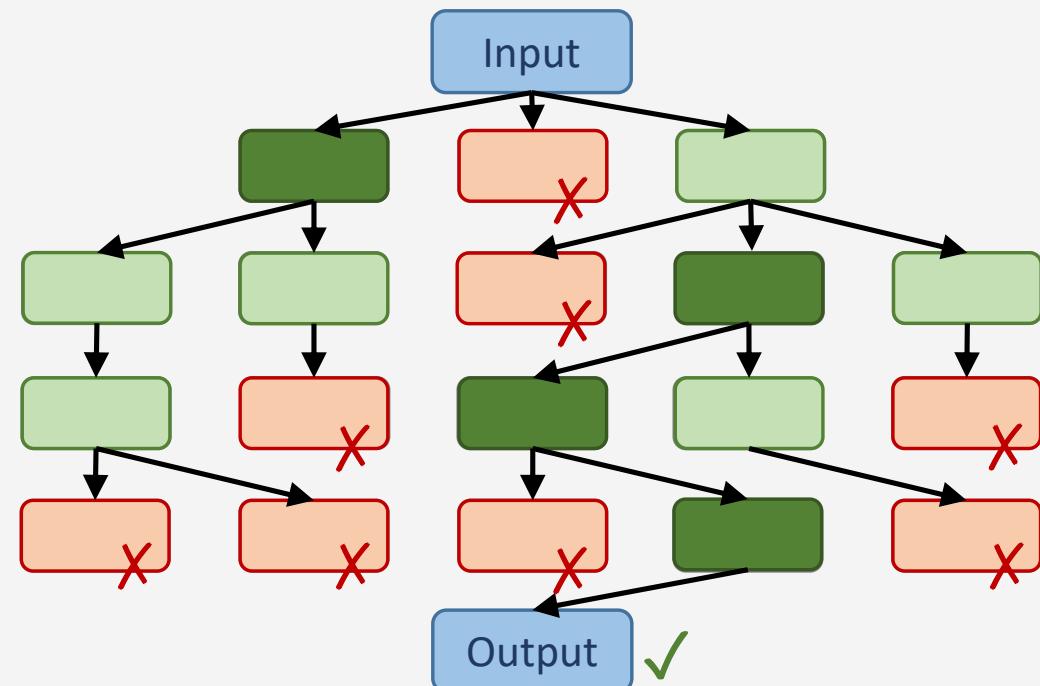
Algorithm 1 ToT-BFS($x, p_\theta, G, k, V, T, b$)

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

```

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

It is an excellent choice when the depth of the tree is relatively small, and solutions are spread out evenly



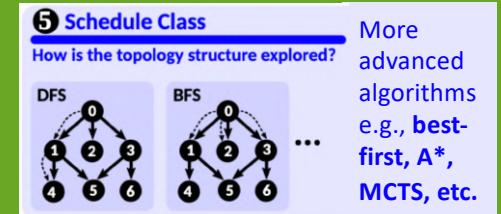
(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>

4)

Search Algorithm: Explore the Solution Space

- Plug and play a proper tree search algorithm based on the tree structure
- **BFS:** maintaining a set of the b most promising states per step
- **DFS:** exploring the most promising state first
 - Until the final output is reached or state evaluator find the thought impossible to succeed



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

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

if $t > T$ **then** record output $G(p_\theta, s, 1)$

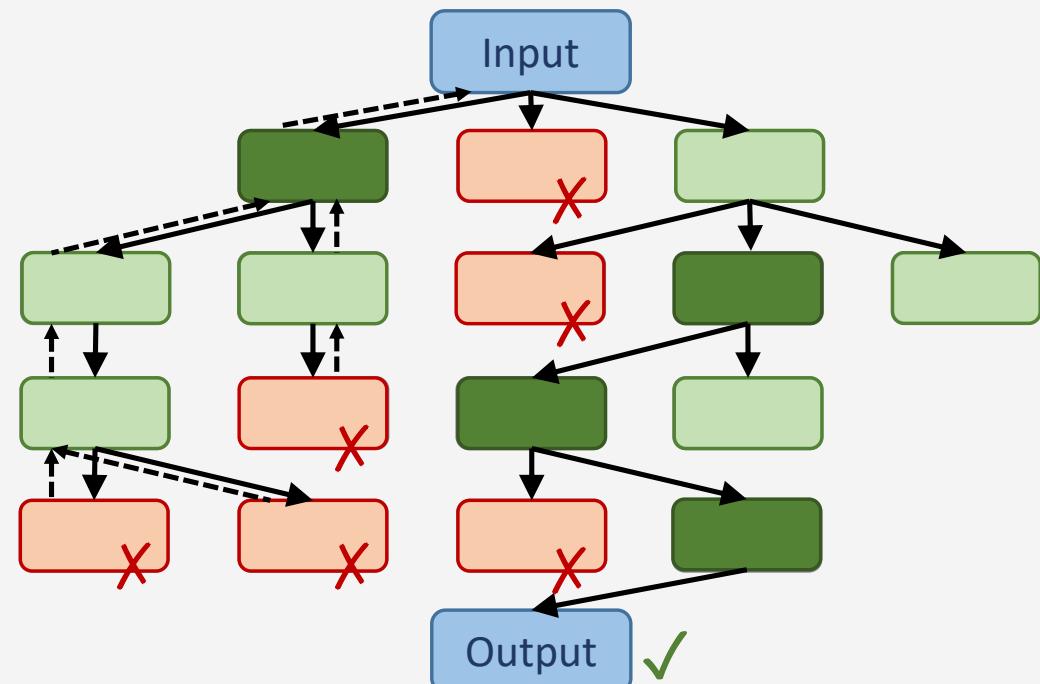
end if

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

end if

end for

It is suitable when the tree depth is significant, and solutions are located deep in the tree



(Yao et al. 2023) <https://arxiv.org/abs/2305.10601>
<https://arxiv.org/html/2312.17445v2>
<https://cameronrwolfe.substack.com/p/tree-of-thoughts-prompting>

<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>

```

# Initialize an agent from swarms
agent = Agent(
    agent_name="tree_of_thoughts",
    agent_description="This agent uses the tree_of_thoughts library to generate thoughts.",
    system_prompt=None,
    llm = OpenAIChat(),
)

# Initialize the ToTAgent class with the API key
model = ToTAgent(
    agent,
    strategy="cot",
    evaluation_strategy="value",
    enable_react=True,
    k=3,
)

# Initialize the MonteCarloSearch class with the model
tree_of_thoughts = MonteCarloSearch(model)

# Define the initial prompt
initial_prompt = ""

Input: 2 8 8 14
Possible next steps:
2 + 8 = 10 (left: 8 10 14)
8 / 2 = 4 (left: 4 8 14)
14 + 2 = 16 (left: 8 8 16)
2 * 8 = 16 (left: 8 14 16)
8 - 2 = 6 (left: 6 8 14)
14 - 8 = 6 (left: 2 6 8)
14 / 2 = 7 (left: 7 8 8)
14 - 2 = 12 (left: 8 8 12)

Input: use 4 numbers and basic arithmetic operations (+-*/) to obtain 24 in 1 equation
Possible next steps:
.....

```

Generator & Evaluator

Search Algorithm

Few-shot CoT Propose Prompt

Tree-of-Thought Piecing Together

```

# Define the number of thoughts to generate
num_thoughts = 1
max_steps = 3
max_states = 4
pruning_threshold = 0.5

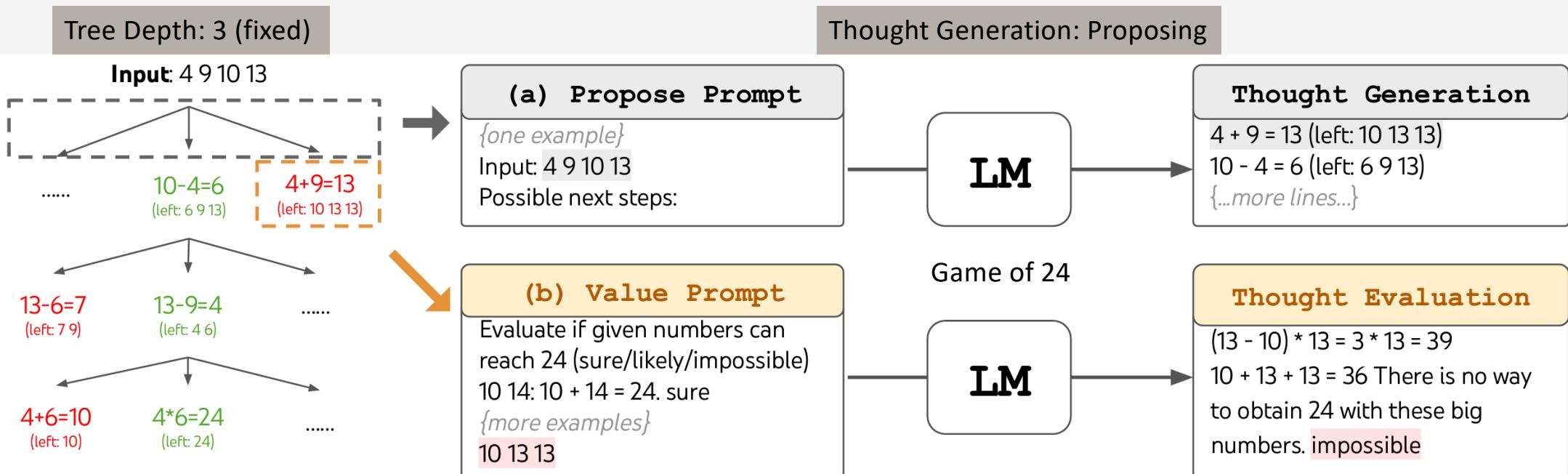
```

+
Generate the thoughts
solution = tree_of_thoughts.solve(
 initial_prompt=initial_prompt,
 num_thoughts=num_thoughts,
 max_steps=max_steps,
 max_states=max_states,
 pruning_threshold=pruning_threshold,
 # sleep_time=sleep_time
)

print(f"Solution: {solution}")

https://github.com/kyegomez/tree-of-thoughts?source=post_page----f2f744786e65

Tree-of-Thought Prompting Example I: Game of 24



Thought: intermediate equations and the corresponding outcome

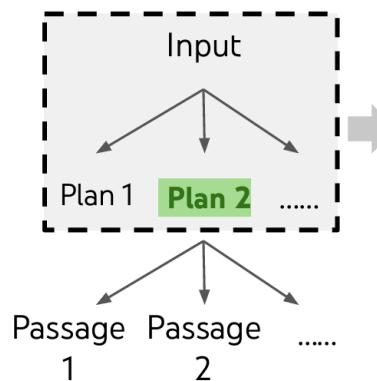
State Evaluation: Value (Classification)

Search Algorithm: BFS



Tree-of-Thought Prompting Example II: Creative Writing

Tree Depth: 1 (fixed, if plan only); 2~ (varied, if including content)



(a)
Input

(b)
Plans

(c)
Votes

Thought Generation: Sampling

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.

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

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.**

Thought: writing plan (and passages if including content)

State Evaluation: Vote

Search Algorithm: N/A (fixed, if plan only); DFS (varied, if including content)



Support to Problem Solving

- **Contextual depth:** providing clear direction
- **Structured exploration:** facilitating systematic analysis
- **Enhanced relevance:** tailoring responses to angles
- **Iterative refinement:** optimizing for accuracy and coherence

Other Possible Variations

- **Typology derivation:** manually defined vs. constructed by LLM
- **Typology representation:** implicit vs. explicit
- **Search representation:** textual description, algorithmic code, in-context example(s), etc.

- Cost
- Complexity
- Information Load
- Linear Relationship

https://www.youtube.com/watch?v=n10_cdSGjXY

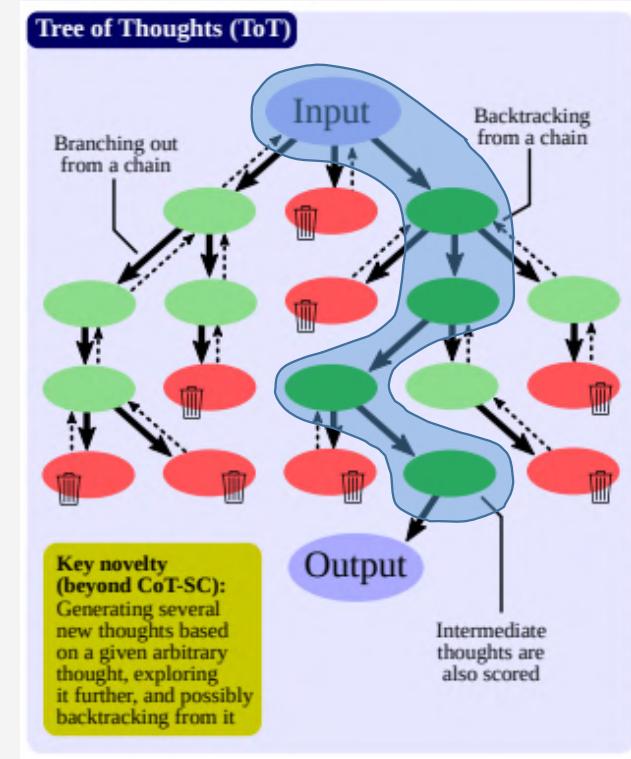
<https://www.semanticscholar.org/reader/eff9d7ed06f30f121d30ee13802a11f172ef66f4>

(Besta et al., 2023) <https://arxiv.org/pdf/2401.14295.pdf>

Assuming a **fixed thought size** (#tokens) and a **fixed context size N** (#thoughts in the LLM context)

Latency: Number of steps between the input and output

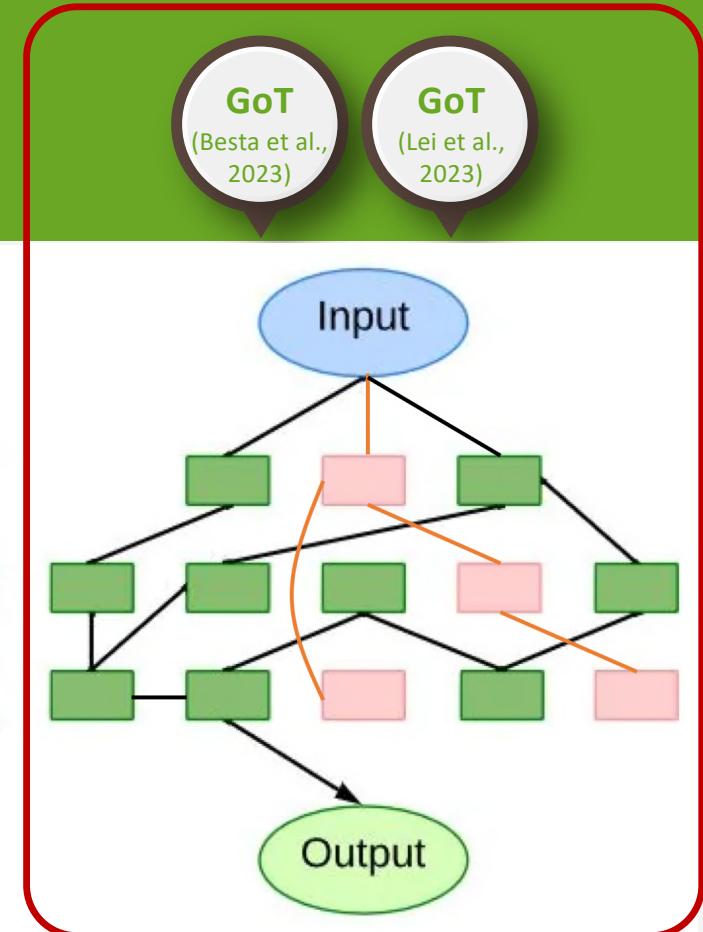
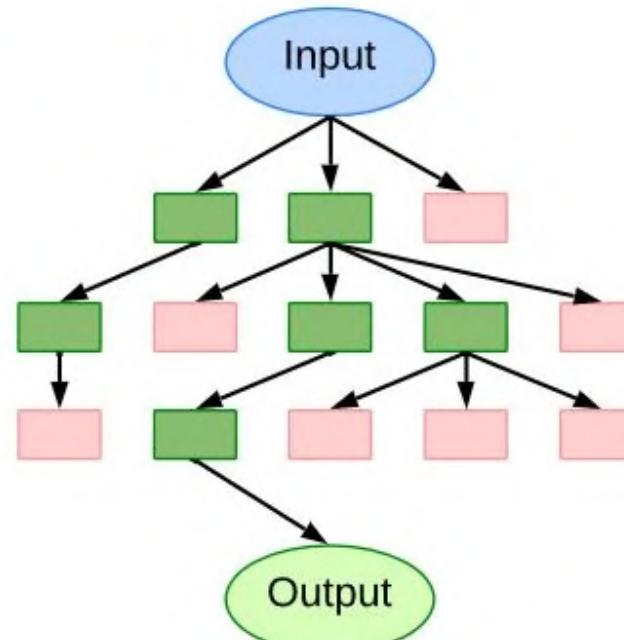
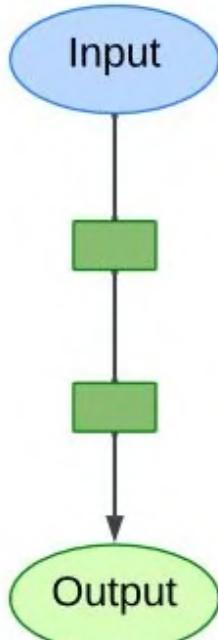
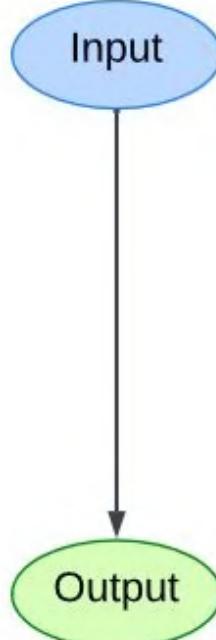
Volume: For a thought t , the number of preceding thoughts that could impact t



XX Low Volume: $O(\log_k N)$

✓✓ Low Latency: $O(\log_k N)$

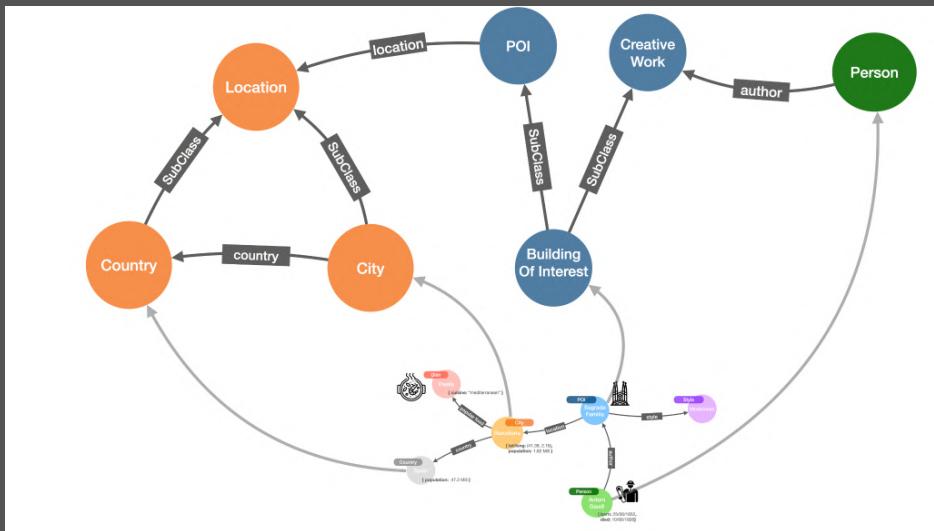
Prompt Engineering Techniques (Prompting Paradigms)



04

Simulate Human Thinking Graph-of-Thoughts

- Description of how a graph-of-thoughts works
- Concepts and mechanisms of graphs-of-thoughts
- Pros and cons of using graphs-of-thoughts



“...model the LLM reasoning as an arbitrary graph, where thoughts are vertices and dependencies between thoughts are edges” (Besta et al., 2023)

04

Simulate Human Thinking

Graph-of-Thoughts

Prompter

Prepare the prompt to be sent to the LLM

Parser

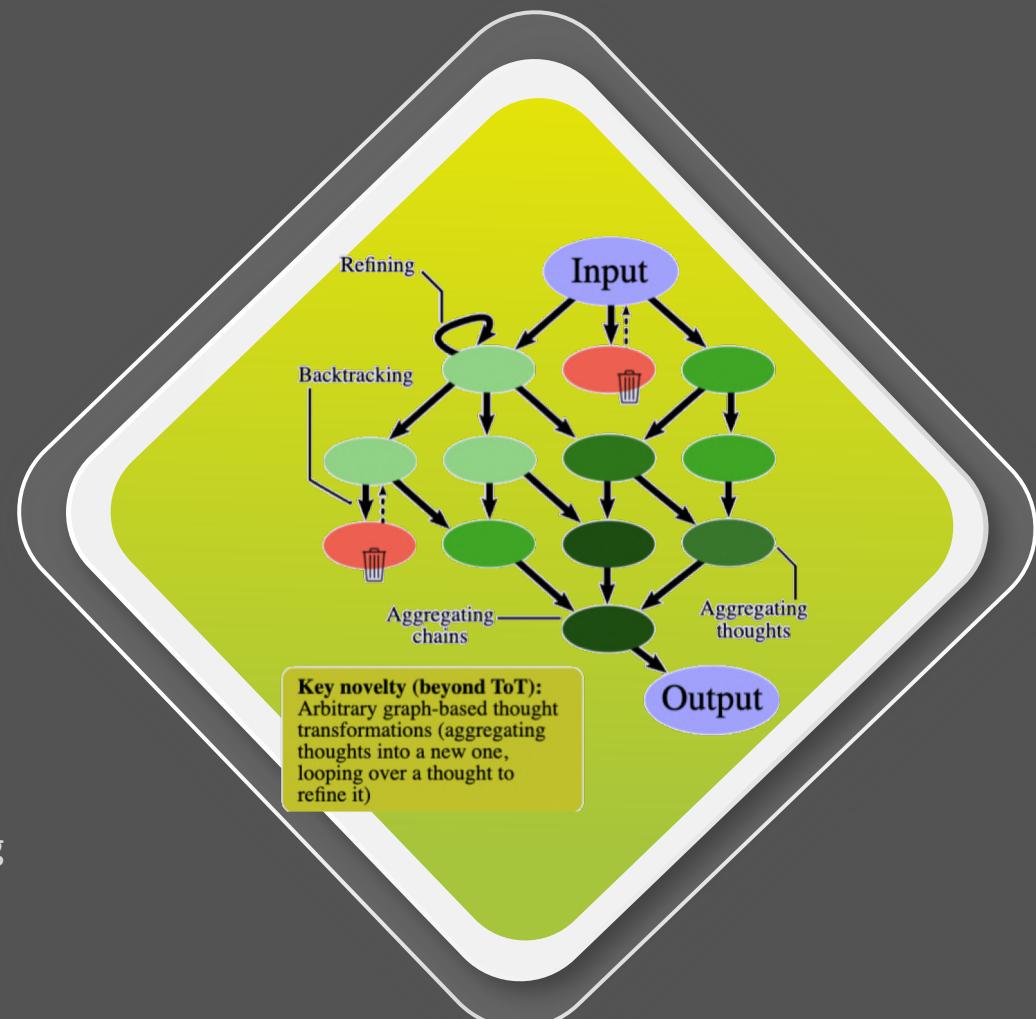
Extract information from the LLM's generated thoughts

Validator

Verify the correctness conditions of a given LLM's thought

Controller

Implements a strategy for selecting thoughts from its Graph of Reasoning Structure (GRS)



Graph-of-Thoughts: Graph-enabled “Transformations of Thought”

“Operation”:

Construction of thought with the previous thought as direct input

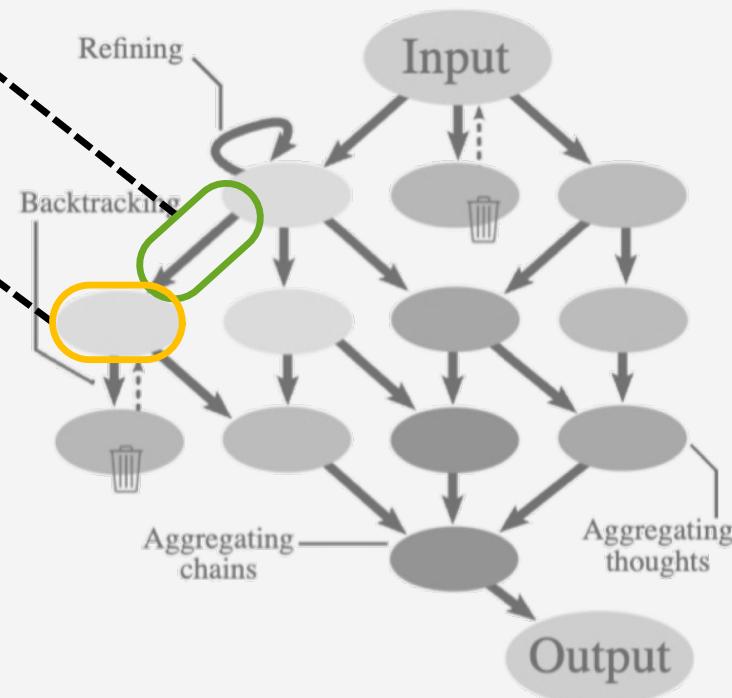
“Thought”:

Solution to a problem at-hand

Characteristics of GoT

- Directed Graph
- Can be heterogeneous

E.g., In creative writing, some vertices can be the writing plans while others can be the actual paragraphs of text

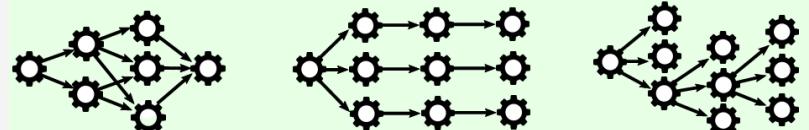


“Graph of Operations (GoO)”:

A static structure that specifies the graph decomposition of a given task

Specifying the Structure of Graph of Operations (GoO)

Graph of Operations enables seamless specification of not only GoT, but also existing schemes such as CoT, CoT-SC, ToT



(Besta et al. 2023) <https://arxiv.org/pdf/2308.09687v2>

Graph-of-Thoughts: Graph-enabled “Transformations of Thought”

“Operation”:

Construction of thought with the previous thought as direct input

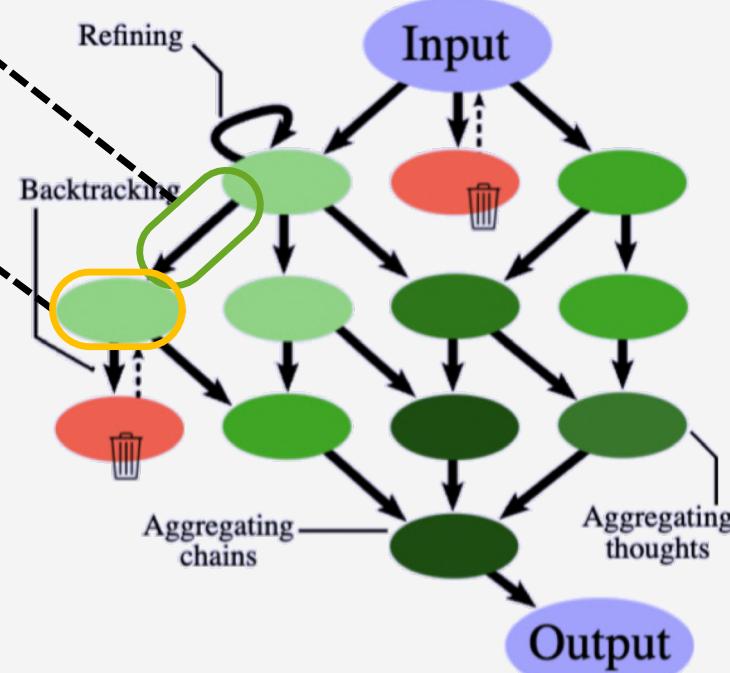
“Thought”:

Solution to a problem at-hand

Characteristics of GoT

- Directed Graph
- Can be heterogeneous

E.g., In creative writing, some vertices can be the writing plans while others can be the actual paragraphs of text



(Besta et al. 2023) <https://arxiv.org/pdf/2308.09687v2>

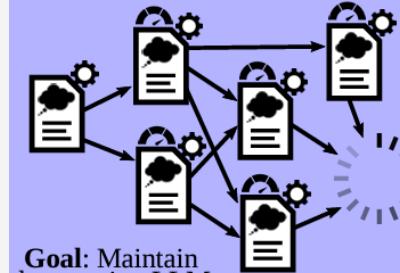
“Graph of Operations (GoO)”:

A static structure that specifies the graph decomposition of a given task

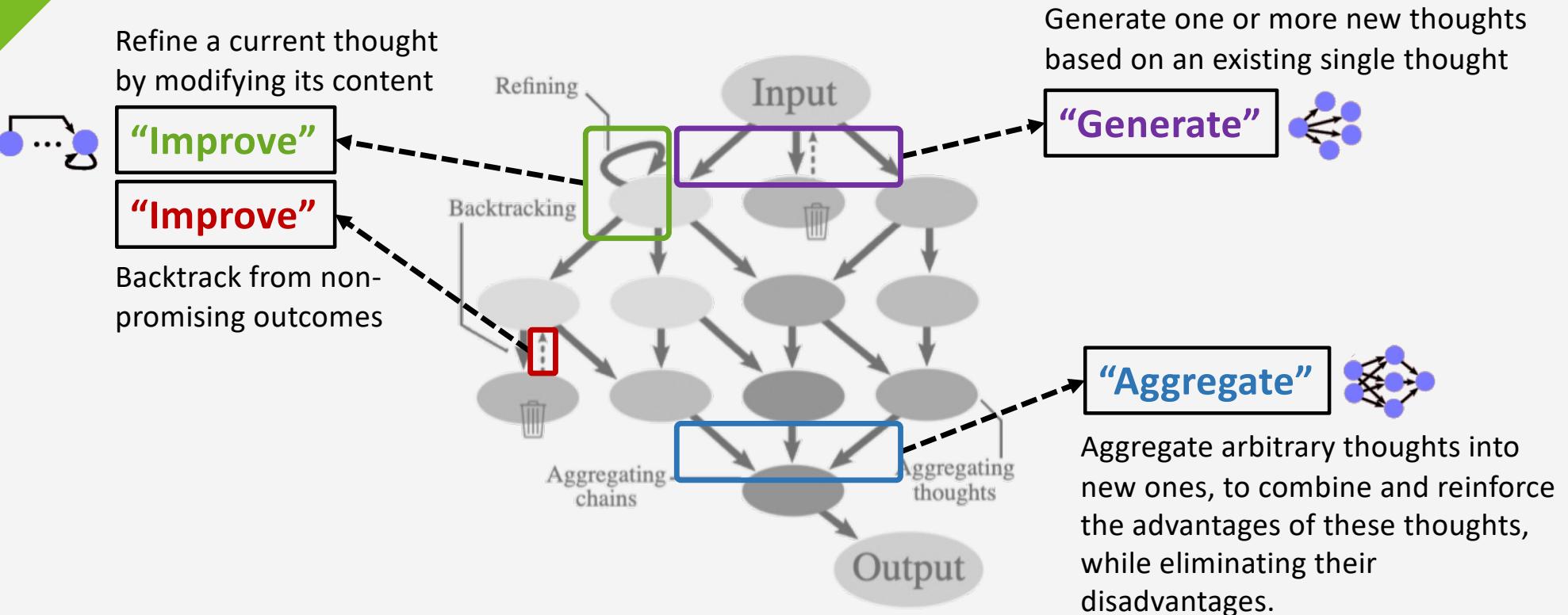
“Graph Reasoning State (GRS)”:

A dynamic structure that maintains the current state of reasoning, e.g., states of individual thoughts (i.e., the intermediate outcomes of operations), their validity, and scores

Graph Reasoning State



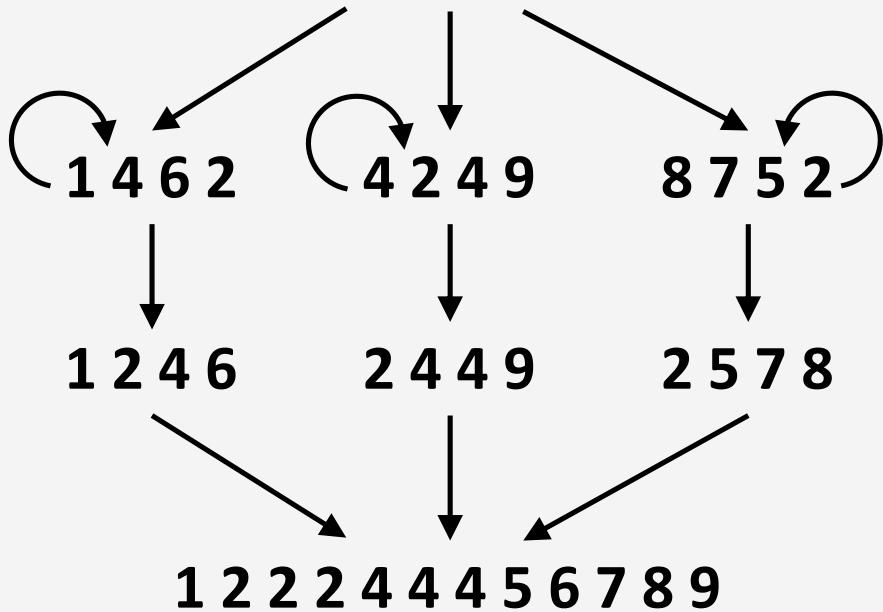
Graph-of-Thoughts: Graph-enabled “Transformations of Thought”



(Besta et al. 2023) <https://arxiv.org/pdf/2308.09687v2>

Examples of Graph-of-Thoughts

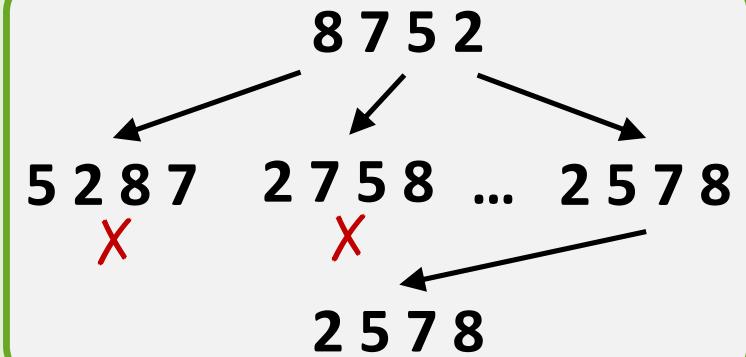
Sort 1 4 6 2 4 2 4 9 8 7 5 2



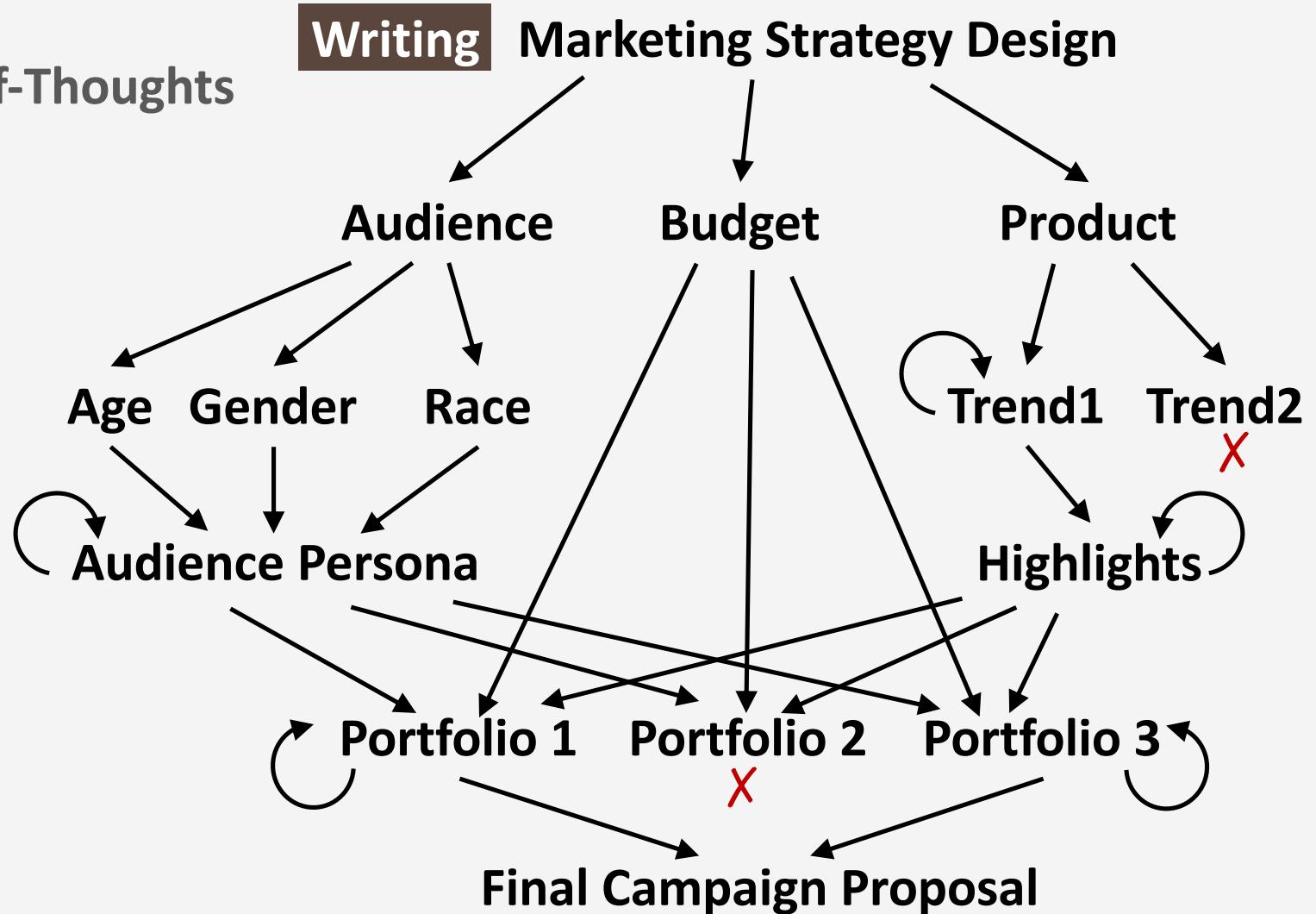
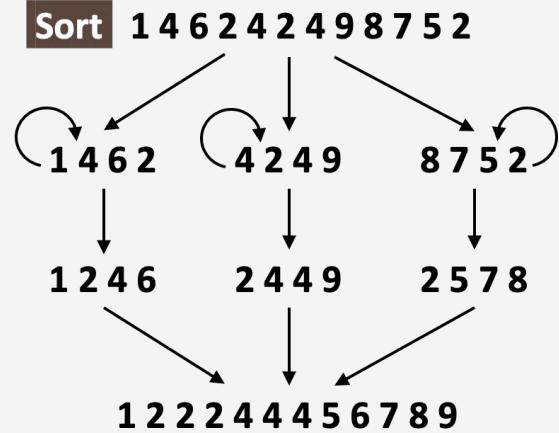
“Generate”

“Improve”

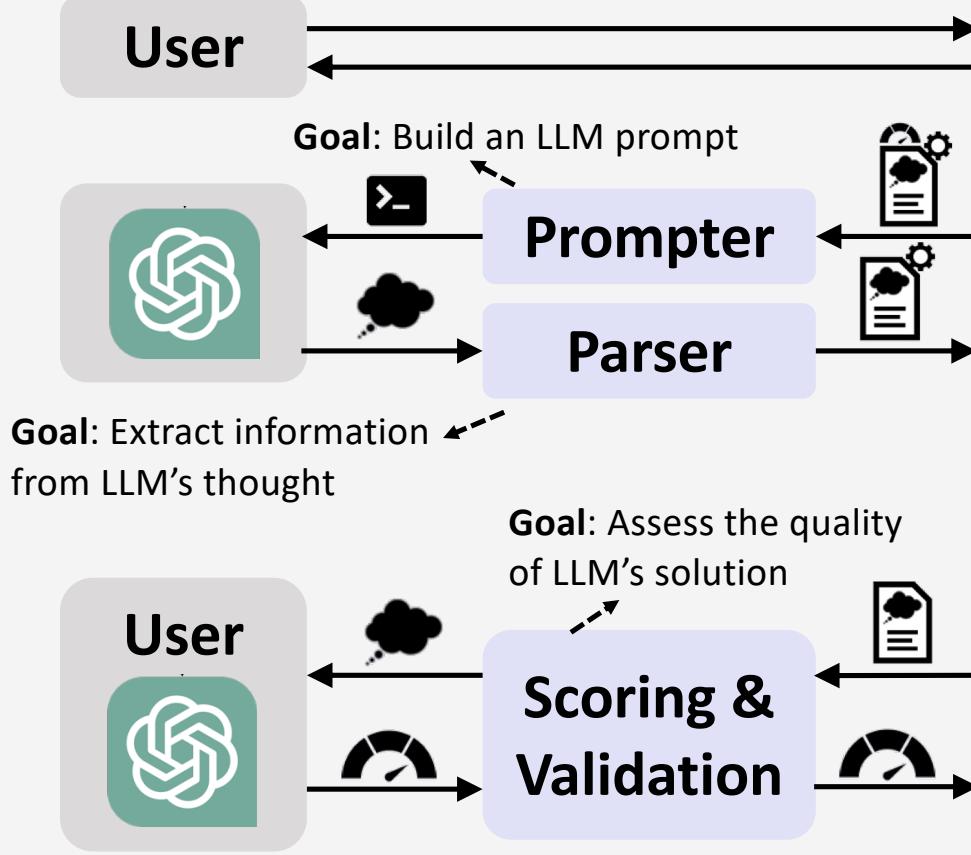
“Aggregate”



Examples of Graph-of-Thoughts



Graph-of-Thoughts Architecture



Goal: Initiate, coordinate, manage, and progress of the GoT execution

Controller

Goal: Specify LLM Thought Transformation

Graph of Operations

Graph Reasoning State

Goal: Maintain the ongoing LLM process

Ranking

Goal: Indicate the top-scoring thoughts

	Thought
	Initial, intermediate, or final solution
	Thought state
	Thought + meta info
	Prompt
	Simple prompt, CoT, ToT, etc.
	Operation
	Generate, Score, Validate, Improve, Repeat, Aggregate, etc.
	Score
	LLM or human

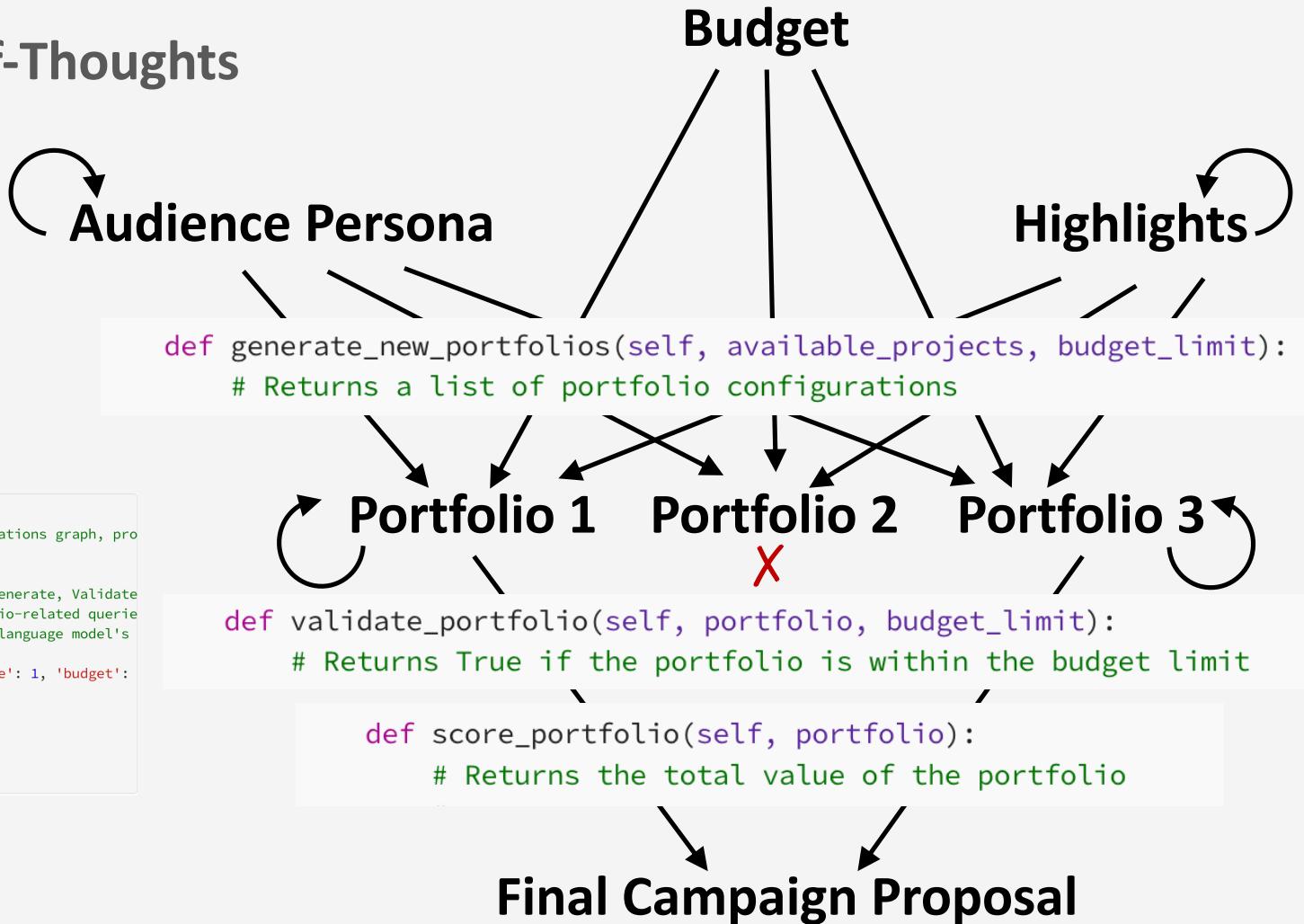


Examples of Graph-of-Thoughts



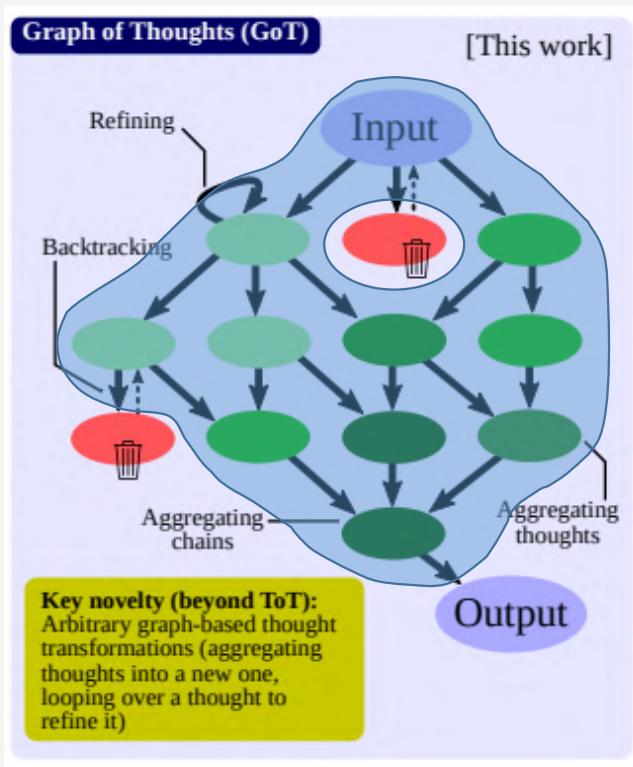
```

# Initialization of the Controller with the language model, operations graph, problem
controller = Controller(
    lm=ChatGPT(), # Instance of the language model
    graph=graph_of_operations, # Operations graph composed of Generate, Validate
    prompter=PortfolioPrompter(), # Custom prompter for portfolio-related queries
    parser=PortfolioParser(), # Custom parser to interpret the language model's
    problem_parameters={
        'available_projects': [{"value": 3, 'budget': 2}, {"value": 1, 'budget': 1},
        'budget_limit': 2,
        'generated_portfolio': None
    }
)
  
```



Assuming a **fixed thought size** (#tokens) and a **fixed context size N** (#thoughts in the LLM context)

Latency: Number of steps between the input and output



✓✓ Large Volume: N

✓✓ Low Latency: $O(\log_k N)$

GoT Key Characteristics

- Suitable for elaborate problem cases
- Higher quality of outcome
- Higher cost

GoT Areas of Improvement

- Balance decomposition of tasks
 - Size of graph
 - Executability
 - Accuracy to avoid improve operation
- Reduce static prompt overhead
 - Difficulty in defining the GoO
 - Decrease the size of few-shot examples
 - ...

https://www.youtube.com/watch?v=n10_cdSGjXY



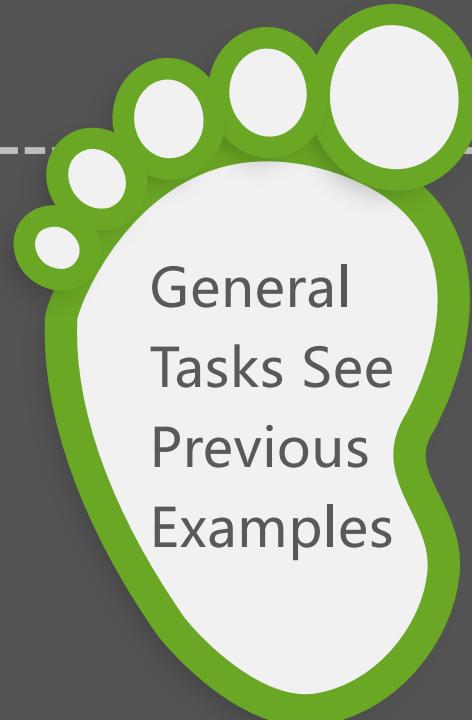
Break

We will be back for exercise in ... minutes

Applications of CoT, ToT, and GoT in Conversational Interactions

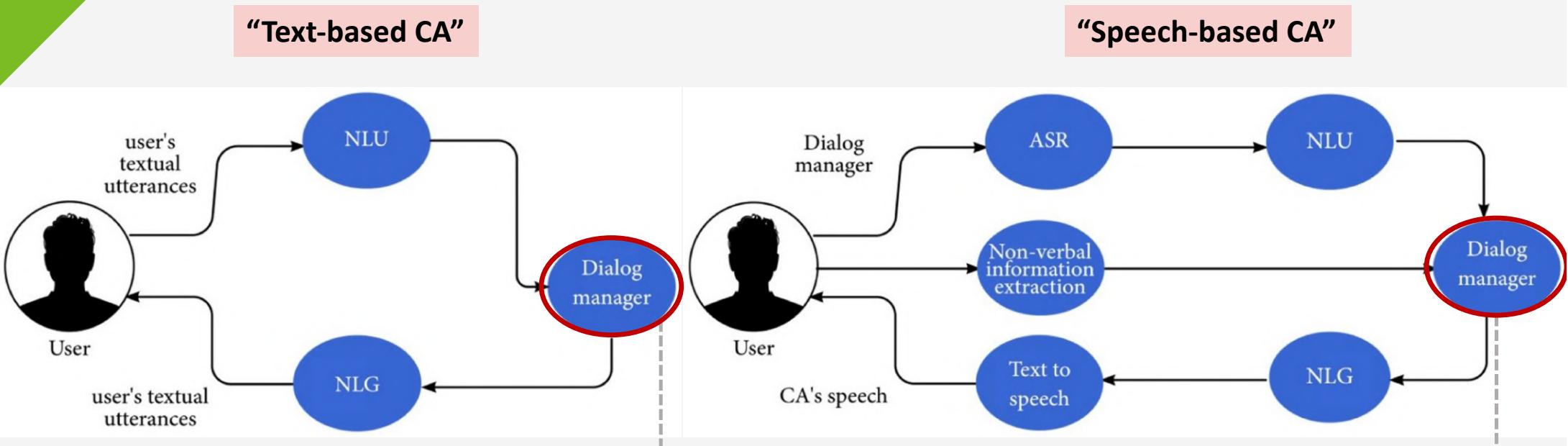
General Task Scenarios

- Arithmetic reasoning
- Common sense reasoning
- Content creation
- Data analysis
- Dialogue task
 - Give instruction
 - Collect info
 - Q&A
 - ...
- ...



- **Dialogue Task Domains**
Conversational
Interactions
- Customer service
 - Education (e.g., tutor)
 - Healthcare (e.g., nurse)
 - ...

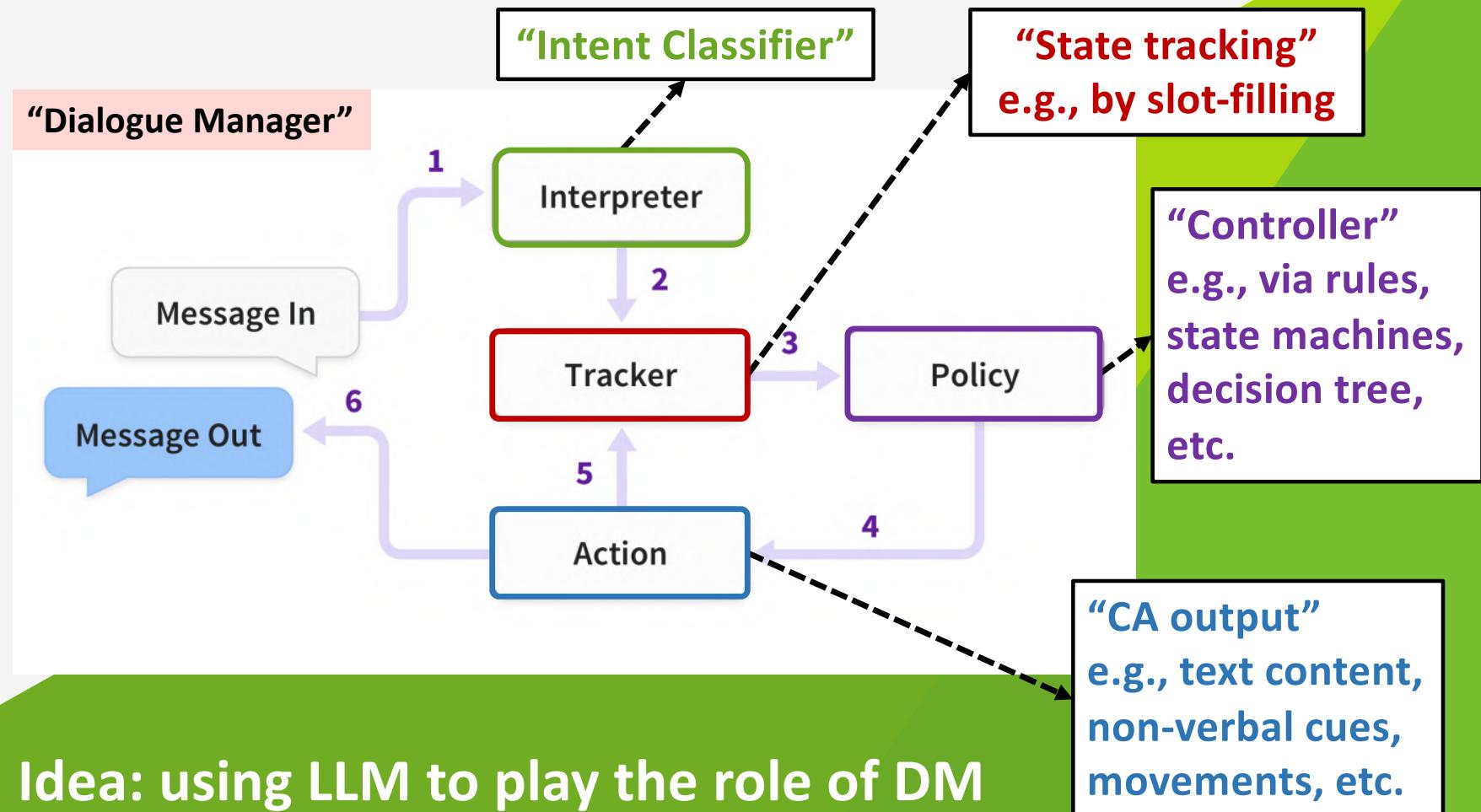
Classic Architecture of a Conversational Agent (CA)



The aim of a Dialogue Manager (DM)
is to mimic all cognitive aspects
related to a natural conversation,
handling dialogue states and flow.

A Close-up Look at a Dialogue Manager

A Dialogue Manager interprets users' intents and dialogue context, keeps track of the current state of the conversation, selects the policy to proceed with, and controls the agent's output



Idea: using LLM to play the role of DM



Exercise



Chain-based schemes introduce explicit intermediate LLM thoughts between the input and the output. This linear sequence of thoughts guides the LLM in a step-by-step manner towards the solution, enhancing the clarity and traceability of the reasoning process.



Graph-based schemes offer an arbitrary reasoning framework. They enable the aggregation of various reasoning steps into a synergistic solution, allowing for non-linear and multifaceted problem-solving approaches.

Tree-based schemes bring the possibility to explore several next-step variants at each juncture, allowing the LLM to evaluate multiple pathways and select the most promising one. This branching structure facilitates a broader exploration of potential solutions.

↑ Complexity
↑ Flexibility
↑ Cost and Load

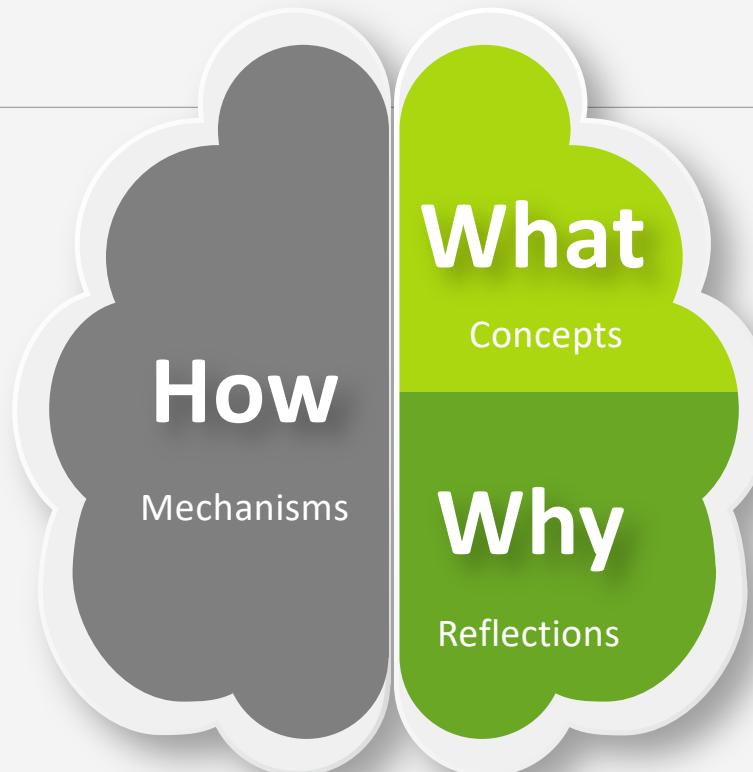


Conclusion



Summary

- Fundamentals
 - Introduction of Chain-of-Thought, Tree-of-Thoughts, and Graph-of-Thoughts
- Examples
 - How to apply them in general and dialogue tasks
- Discussion
 - Pros and cons



Future Directions

- Improve cost-effectiveness,
 - E.g., encoding the structure within a single prompt
- Explore new architectures
 - Productivity
 - Programmability
 - Scalability
 - Parallelizability
- Develop understanding
 - Theoretical
 - Empirical



Charting the Routes of Thoughts in LLM-Empowered Conversational Interactions

Xiaojuan Ma

Human-Computer Interaction Initiative
Hong Kong University of Science and Technology
mxj@cse.ust.hk

Special Thanks: Dingdong Liu

CIX'24