

# Why Prompt Design Matters and Works: A Complexity Analysis of Prompt Search Space in LLMs

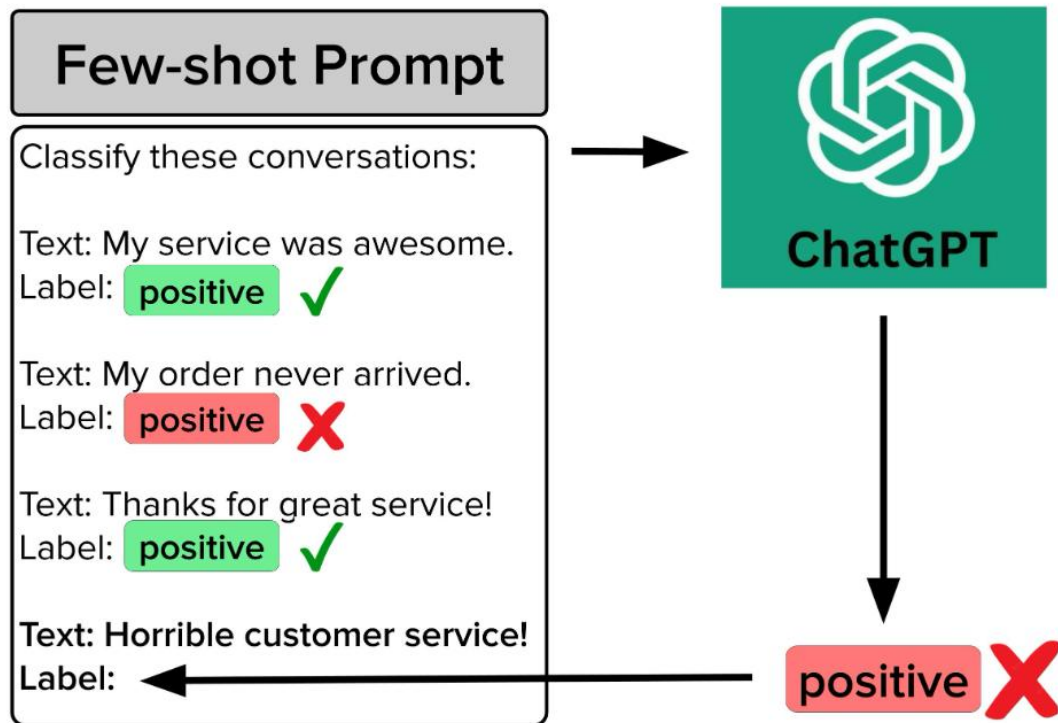
---

**Xiang Zhang<sup>1\*</sup>**    **Juntai Cao<sup>1\*</sup>**    **Jiaqi Wei<sup>3</sup>**    **Chenyu You<sup>2†</sup>**    **Dujian Ding<sup>1†</sup>**

<sup>1</sup>University of British Columbia, <sup>2</sup>Stony Brook University, <sup>3</sup>Zhejiang University  
xzhang23@ualberta.ca, {jtcao7, dujian}@cs.ubc.ca,  
jiaqi.wei@zju.edu.cn, cyou@cs.stonybrook.edu

**ACL 2025**

**Prompt:** An instruction, question, or input text that a user gives to an AI model to clearly tell the model what task needs to be completed and how to output the results.



## Chain-of-Thought Prompting

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

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

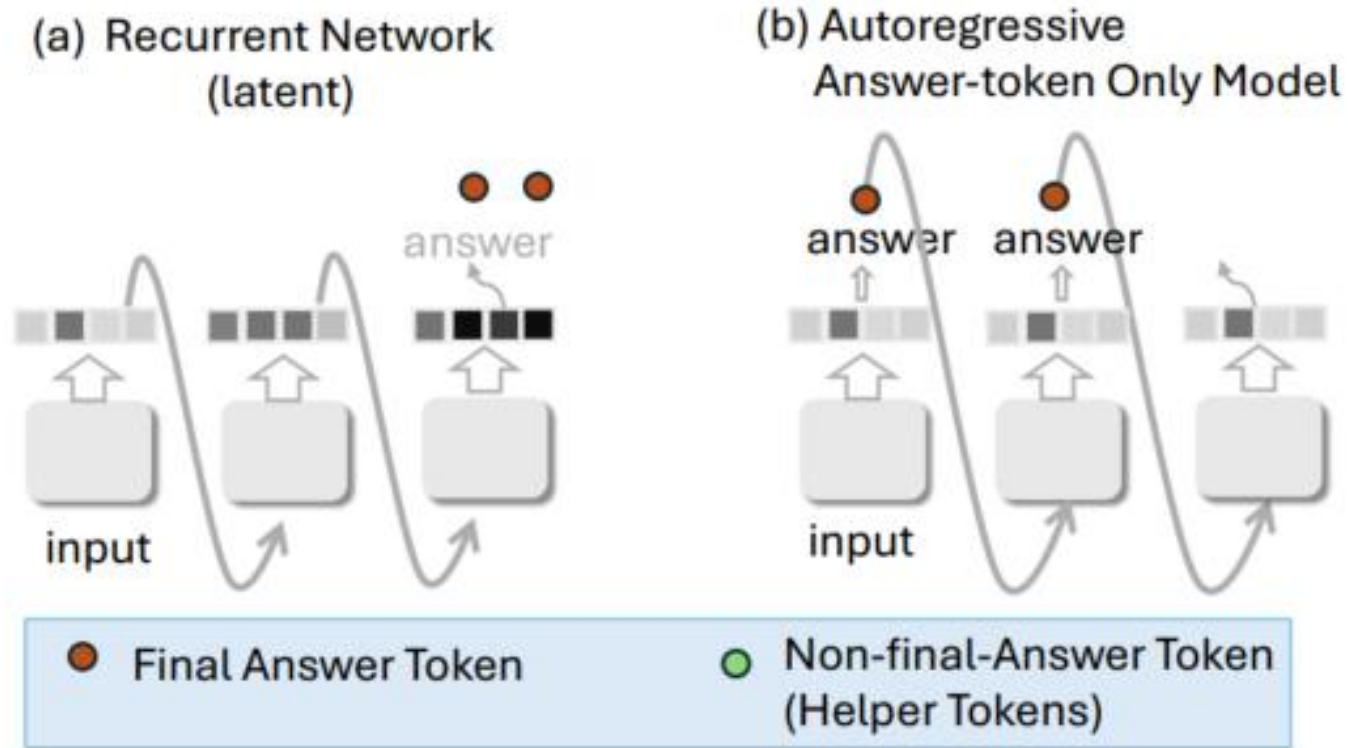
Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

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

- Why does a particular prompt design work?
- How to find the optimal prompt design?

# Limitations of Transformer Architecture



$h$  stores reasoning memory and intermediate reasoning results  
the ability to sequentially compute and update  $h$  over time  
allows a model to build reasoning depth

Answer  $y$  is a discrete value extracted from  $h$  and only  
contains partial information

$y$  exists outside the latent space where  $h$  operates, meaning it  
cannot be used for computation

**RNN:** Sequential token processing

$$h_t = g_\theta(h_{t-1}, x_t)$$

This kind of 'recursive transmission' can  
naturally accumulate temporal information

**Transformer:** Self-Attention

$$h_t = \text{Attention-weighted sum of all tokens}$$

**For RNN:** The longer the sequence, the more  
times the hidden state is recursively updated,  
and the inference depth naturally increases.

**For Transformer:** The computation of its hidden  
state depends only on the 'number of layers' of  
the model and is independent of the length of  
the input sequence.

# Nature of Inductive Reasoning



$$x_n = (x_1, x_2, \dots, x_n)$$

To validate the  $n$ -th move, the  $n$ -th board state  $h_n$  must be calculated.

This requires  $n$  sequential computations, as the  $n$ -th board state depends not only on the sequence of moves  $x$  but also on the previous board state  $h_{n-1}$ .

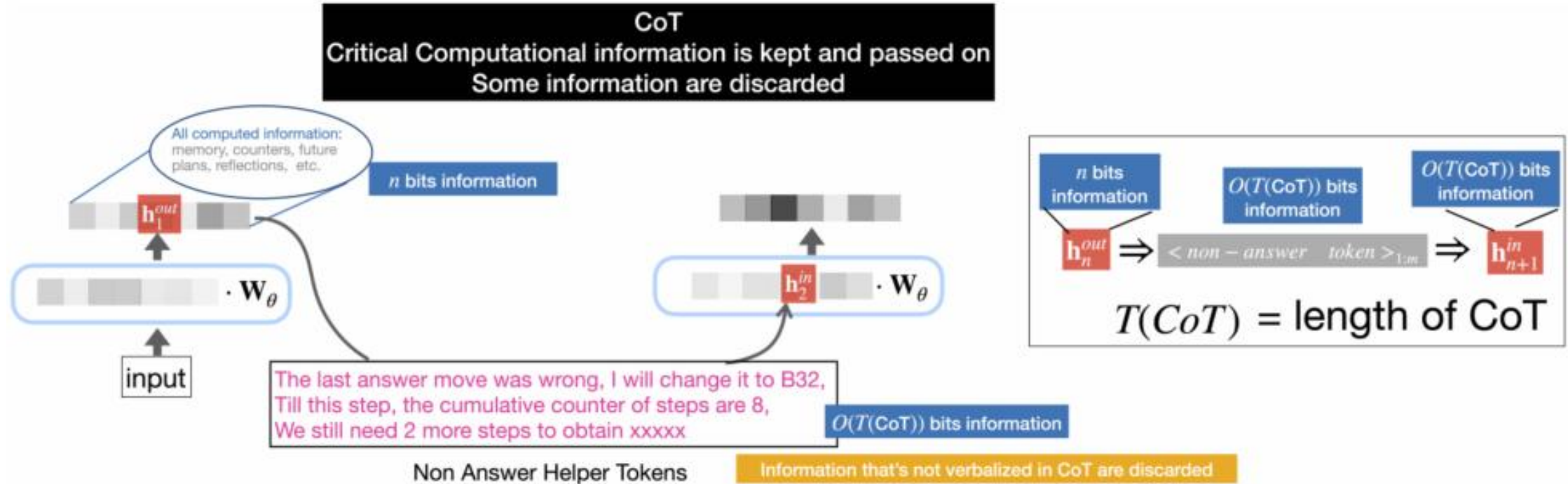
While a neural network could memorize the mapping from  $x_{1:n}$  to the correct  $h$ , the memorization will be an exponentially growing challenge

## Reasoning inherently requires sequential depth

To solve a given task, there is a theoretical lower bound on the required depth of computation

Transformer's fixed sequential reasoning depth over the hidden state  $h$  prevents them from solving tasks that require deeper reasoning as input length grows

# CoT+Autoregressive=Recurrent



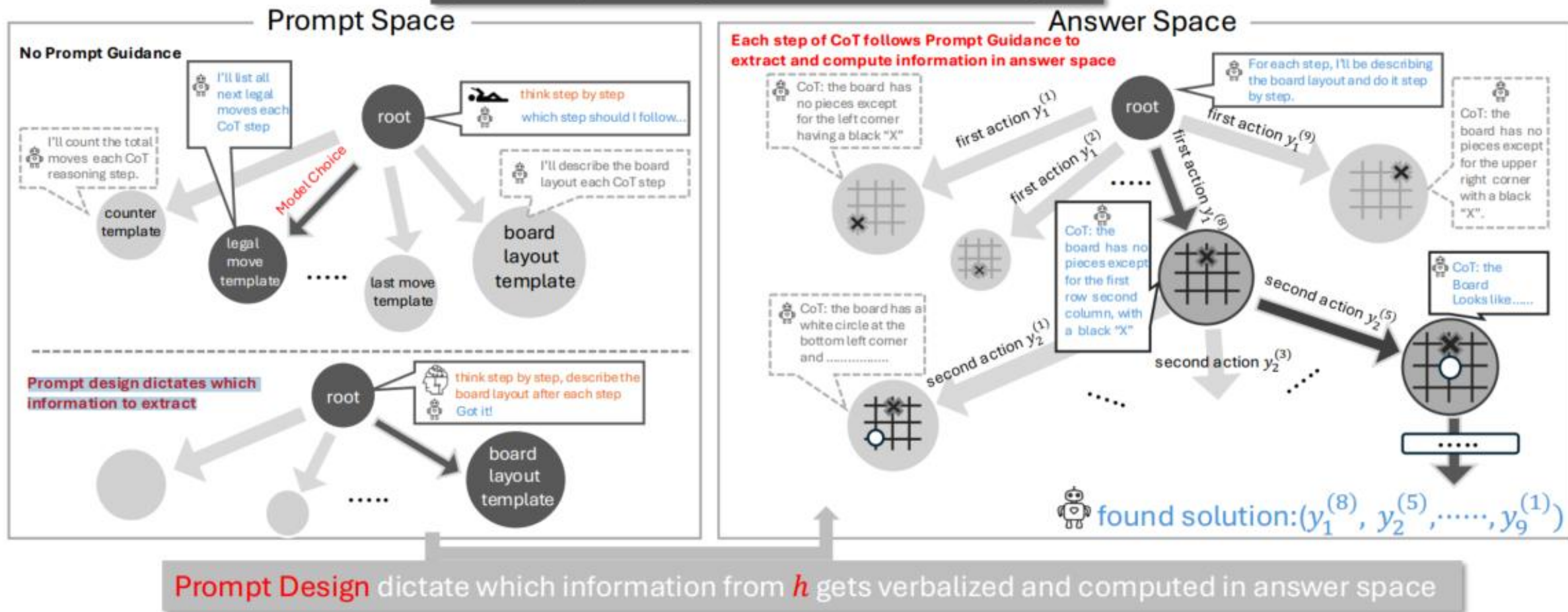
CoT extends beyond simple answer token generation by producing intermediate steps as non-answer natural language tokens  $(o_1, o_2, \dots, o_k)$  that act as a discretizations of latent information  $h_n$

$$h_t \xrightarrow{\text{discretization}} (o_1, o_2, \dots, o_k) \xrightarrow{\text{vectorization}} h_{t+1}$$

LLMs with CoT extend reasoning from latent space  $H$  to natural language token space  $O$ .

# CoT Search Space = Prompt Space + Answer Space

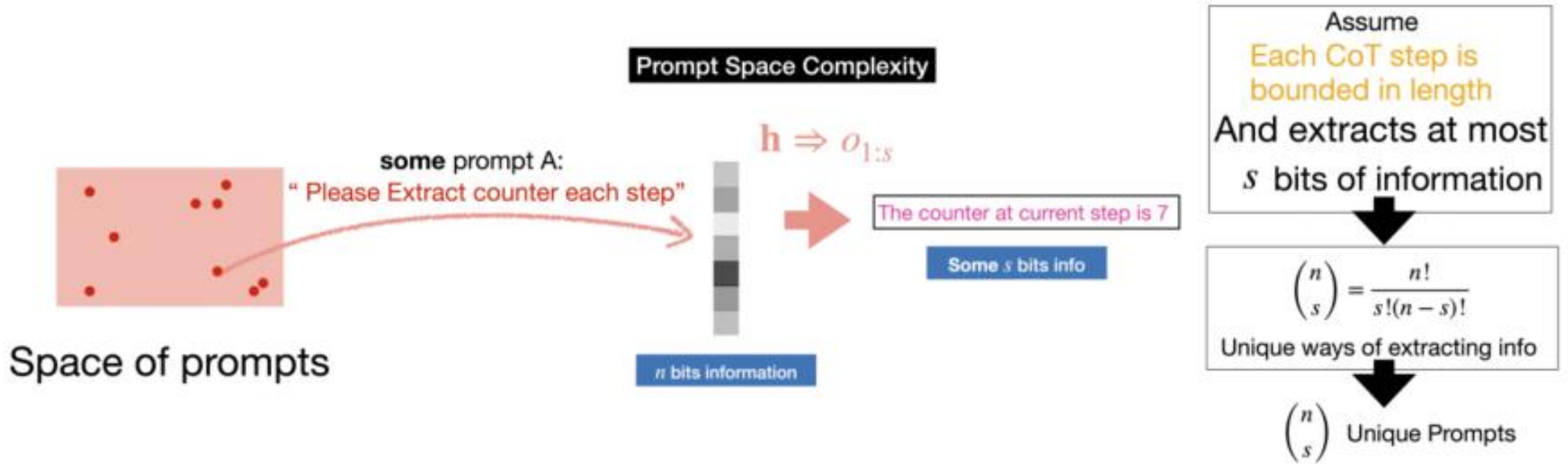
QA: find a sequence of legal moves that leads to the end game



**Decompose the CoT reasoning into two components:**

1. template search within the prompt space
2. answer search within the answer space

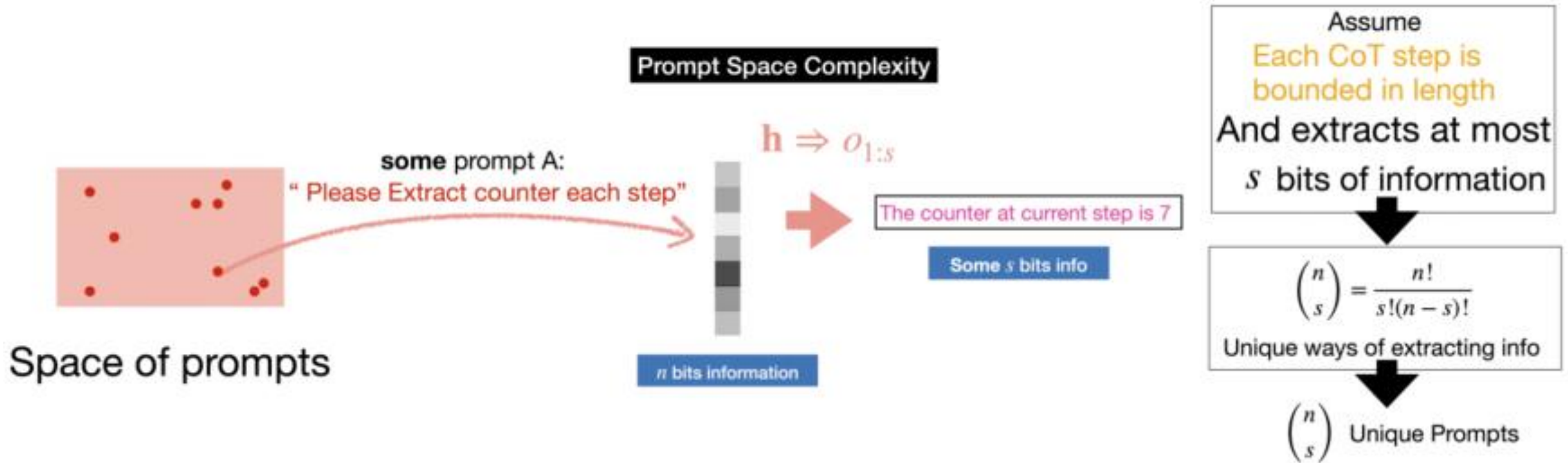
# Prompt Space Complexity



When LLMs are prompted to perform tasks, they follow a step template, specifying which information from  $h$  to extract and discretize into non-answer helper tokens  $(o_1, o_2, \dots, o_{T(\text{CoT})})$  in CoT.

Ideally, as  $T_{(\text{CoT})} \rightarrow \infty$ —meaning the length of the each CoT step is arbitrarily long—all vectorized information in  $h$  can be fully textualized

# Prompt Space Complexity



Define the amount of information stored in  $h$  as  $n$  bits, and each CoT step extracts up to  $s$  bits of information into  $o$ , each unique step template specifies a way to extract  $s$  bits from the full  $n$ -bit space

$$C(n, s) = \frac{n!}{s!(n-s)!}$$

The prompt template defines how information is extracted and used recurrently in the CoT process. Finding the correct template is equivalent to discovering the algorithm for solving a given task

With a specific step (prompt) template  $p_i$  chosen from the prompt space  $\mathbf{P}$ , CoT iteratively executes

$h_t \xrightarrow{p_i} (o_1^{(i)}, o_2^{(i)}, \dots, o_k^{(i)}) \Rightarrow h_{t+1}$  to update  $\mathbf{h}$  and calculate the next state, continuing this process until reaching the final state (solution).

**In the chess simulation task of <finding a set of actions leading to game end>:**

the answer space  $\mathbf{S} = (s_1, s_2, \dots, s_\infty)$  contains all possible combinations of action sequences  $\mathbf{s}$

The solution set  $\mathbf{CR} \subset \mathbf{S}$  includes all valid action sequences that lead to the end of the game, being a subset of the entire answer space  $\mathbf{S}$

Solving the problem requires identifying one single correct action sequence  $s_{correct} = (y_1, y_2, \dots, y_T) \in \mathbf{CR}$

The complexity of navigating the answer space can be roughly measured by

$$\frac{\text{len}(\mathbf{CR})}{\text{len}(\mathbf{S})} | p$$

**The size of answer space  $\mathbf{S}$ :** the more open the task, the more possible outcomes, and the larger  $\mathbf{S}$  is.

**The size of correct solution set  $\mathbf{CR}$ :** the stricter the task, the fewer correct solutions there are, and the smaller the  $\mathbf{CR}$ .

The prompt template ( $p$ ) directly determines which hidden state information the model 'extracts,' thereby reshaping the complexity of the answer space  $\mathbf{S}$ .

Level	Task	RNN	Tape RNN	Transformer	LLM w/o CoT	CoT Unsupervised	CoT Supervised	CoT Supervised-SUB
<b>R</b>	Modular Arithmetic	<b>1.00</b>	<b>1.00</b>	<b>0.96</b>	0.22	<b>0.96</b>	<b>1.00</b>	0.44
	Parity Check	<b>1.00</b>	<b>1.00</b>	0.52	0.58	<b>0.94</b>	<b>1.00</b>	0.42
	Cycle Navigation	<b>1.00</b>	<b>1.00</b>	0.62	0.50	0.78	<b>1.00</b>	0.26
<b>CF</b>	Stack Manipulation	0.56	1.00	0.58	0.00	<b>0.92</b>	<b>0.96</b>	0.00
	Reverse List	0.62	<b>1.00</b>	0.62	0.48	0.80	<b>0.96</b>	0.38
	Modular Arithmetic	0.41	0.95	0.32	0.00	0.82	<b>0.94</b>	0.50

Solving these tasks demands a minimum computational depth that scales linearly with input length, surpassing the constant depth inherent to Transformer models.

**R** represents simple recursive tasks, **CF** represents complex recursive tasks

**RNN/Tape-RNN** represents models with inherent recursive capabilities

**Transformer** represents a pure Transformer model without CoT assistance

**CoT (Unsupervised)** simply means 'Let's think step by step'

**CoT (Supervised)** refers to the optimal prompt template, directly defined the core algorithm of the task, allowing the model to focus directly on effective reasoning paths without independent exploration.

**CoT (Supervised-SUB)** represents a suboptimal prompt template, extracting irrelevant information or designing redundant steps, causing the model to be unable to effectively propagate intermediate states, and even being misled into incorrect reasoning paths

Model	R			CF		
	MA	PC	CN	SM	RL	MA
Unsupervised <b>CoT</b>	<b>0.96</b>	<b>0.94</b>	0.78	0.92	0.80	0.82
Unsupervised <b>ToT</b>	<b>0.92</b>	<b>0.90</b>	<b>0.92</b>	0.36	0.88	0.78
Unsupervised <b>GoT</b>	<b>1.00</b>	<b>0.98</b>	<b>0.90</b>	0.72	<b>0.92</b>	0.88
Correctly supervised <b>CoT</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.96</b>	<b>0.96</b>	<b>0.94</b>

Q: Under the premise of a fixed prompt template, can different CoT variants improve performance?

The design goal of ToT/GoT is 'to optimize the search method of the answer space based on existing reasoning templates'

But they cannot change the fundamental problem of whether the 'template is effective' — if the template autonomously generated by the model is suboptimal (such as misjudging the push sequence in stack operations), even increasing multi-path exploration will only waste computational resources in the wrong direction of reasoning, and may even reduce performance.

A: It can be slightly optimized, but cannot replace the core role of the optimal prompt template.

**Thanks**