SUPERVISED CHAIN OF THOUGHT

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have revolutionized natural language processing and hold immense potential for advancing Artificial Intelligence. However, the core architecture of most mainstream LLMs-the Transformer-has inherent limitations in computational depth, rendering them theoretically incapable of solving many reasoning tasks that demand increasingly deep computations. Chain of Thought (CoT) prompting has emerged as a technique to address these architectural limitations, as evidenced by several theoretical studies. It offers a promising approach to solving complex reasoning tasks that were previously beyond the capabilities of these models. Despite its successes, CoT and its variants (such as Tree of Thought, Graph of Thought, etc.) rely on a "one-prompt-for-all" approach, using a single prompt structure (e.g., "think step by step") for a wide range of tasks-from counting and sorting to solving mathematical and algorithmic problems. This approach poses significant challenges for models to generate the correct reasoning steps, as the model must navigate through a vast prompt template space to find the appropriate template for each task. In this work, we build upon previous theoretical analyses of CoT to demonstrate how the one-promptfor-all approach can negatively affect the computability of LLMs. We partition the solution process into two spaces: the prompt space and the answer space. Our findings show that task-specific supervision is essential for navigating the prompt space accurately and achieving optimal performance. Through experiments with state-of-the-art LLMs, we reveal a gap in reasoning performance when supervision is applied versus when it is not. Our goal is to provide deeper insights into the mechanisms underlying CoT, offering guidance for the effective design of CoT variants. Additionally, we underscore the limitations of traditional "unsupervised" prompting methods, arguing that users of CoT cannot simply "sit back" and rely entirely on the model. Instead, we advocate for task-specific "supervised" CoT, enriched with human knowledge, to enable more effective reasoning in LLMs.

034

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

1 INTRODUCTION

037

038 The advent of large language models (LLMs) (Achiam et al., 2023) has ushered in a new era for nat-039 ural language processing and artificial intelligence (Kojima et al., 2022; Zhao et al., 2023). These 040 models exhibit remarkable capabilities across various domains (Thirunavukarasu et al., 2023; Wei et al., 2022; Valmeekam et al., 2023; Zhang et al., 2023), achieving near-human performance in tasks 041 such as knowledge retrieval and articulation (Chang et al., 2024). However, concerns have been 042 raised regarding their reasoning abilities (Valmeekam et al., 2022; Zhang et al., 2024). These tasks 043 range from fundamental operations like counting, sorting, and multiplication (Dziri et al., 2024), to 044 more complex challenges such as mathematical problem-solving, algorithm design, and coding (Xu 045 et al., 2022; Thirunavukarasu et al., 2023). Previous research has explored several factors contributing to these reasoning deficiencies, including training optimizations (Thorburn & Kruger, 2022), 047 tokenization methods (Singh & Strouse, 2024), and dataset choices (Ye et al., 2024). Among these, 048 the architecture of the model plays a pivotal role in determining its reasoning capabilities (Raghu et al., 2017; Zhang et al., 2024; Delétang et al., 2022). The backbone architecture of most mainstream LLMs—the Transformer (with finite precision) (Vaswani, 2017)—has intrinsic limitations 051 related to computational depth (Li et al., 2024). Specifically, the attention mechanism within Transformers can perform only a fixed number of sequential computational steps (Li et al., 2024; Zhang 052 et al., 2024; Sanford et al., 2024; Dehghani et al., 2018), leading to constant-depth modeling (Li et al., 2024). As a result, when relying solely on the Transformer's internal reasoning, the model's



Figure 1: (a) Without supervision during CoT, the model generates its own step template for recurrent computation. This template can be incorrect, leading to task failure. (b) With human supervision, the task performance under CoT can be properly guided. (c) When CoT is not employed, the model relies solely on its internal reasoning via the Transformer architecture. (d) The Transformer can only perform constant-depth sequential computations. We assume that this Transformer neither 072 *memorizes* the results nor performs bit-level (circuit) reasoning; instead, reasoning occurs at the neuron (hidden state) level.

073 074

067

068

069

071

075

computability is restricted to TC⁰ (Li et al., 2024; Feng et al., 2024), which confines it to solving 076 tasks of limited complexity and length (Figure 1.c-d). 077

078 The emergence of Chain of Thought (CoT) (Wei et al., 2022) has transformed the computational lim-079 itations imposed by architectural design. Specifically, it overcomes the *constant depth* limitation of 080 a model's internal (latent space) reasoning by extending the reasoning process into the "text" space 081 through prompting (Li et al., 2024; Zhang et al., 2024; Feng et al., 2024). As demonstrated by both 082 theoretical (Li et al., 2024; Zhang et al., 2024; Feng et al., 2024) and empirical research (Li et al., 083 2024; Zhang et al., 2024), CoT effectively enhances the reasoning depth of Transformer-based models (with finite precision), achieving "Turing Completeness" under ideal conditions (Li et al., 2024; 084 Zhang et al., 2024). While the theoretical analysis of CoT focuses on the model's upper bound com-085 putational power, which may not always align with real-world performance, a deep understanding of the CoT mechanism—particularly how it breaks architectural constraints—is crucial for designing 087 better prompts that maximize the computational potential of real-world LLMs. This understanding 880 also forms the basis for our analysis of "supervised" CoT and our prompt search space theories. 089 Therefore, our work first revisits the underlying mechanisms of CoT from a computational perspective, integrating insights from prior research (Li et al., 2024; Zhang et al., 2024; Feng et al., 2024) 091 to demystify them in a clear, yet comprehensive manner, enriched by our unique perspective.

092 Although theoretical analysis has proven the *existence* of solutions for (almost) any problem using 093 CoT, based on computability and Turing Completeness theory, the actual *discovery* of those solu-094 tions can be much more challenging. This is akin to how a Turing machine can model solutions 095 for any problem (Boolos et al., 2002) but *finding* the exact Turing machine for a specific NP prob-096 lem could be difficult. These challenges arise from two main factors for LLMs with CoT. First, the 097 model must develop the correct "step-by-step" template, which essentially embodies the algorithm 098 used for solving the problem (Figure 1.a-b). For instance, the "steps" for solving a graph search problem using depth-first search (DFS) differ from those of a breadth-first search (BFS) algorithm. 099 Second, even after the template (algorithm) is established, finding the solution might require ex-100 tensive reasoning and exploration to achieve the optimal outcomes. For example, using the BFS 101 template to locate a target node in a tree involves traversing multiple paths in the search space that 102 can be computationally expensive and error-prone. 103

104 The vanilla design of CoT is "unsupervised", meaning that the model generates its step template 105 without task-specific supervision from humans. Specifically, when prompted to "think step by step", LLMs autonomously generate a step template (algorithm) it needs to follow-for instance, generat-106 ing previously visited paths at each step—and then proceeding to search for answers based on this 107 self-generated template (Figure 1.a). Clearly, this naive CoT approach can lead to poor performance,



Figure 2: Comparison between recurrence and autoregression.

as the model may generate sub-optimal step templates (algorithms), which hinder the search process. For example, a problem requiring DFS might be unnecessarily attempted with a BFS template
 generated by the vanilla CoT, incurring high inference costs and likely delivering incorrect answers
 (Figure 1.a).

126 Variants of Chain of Thought, such as Tree-of-Thought (Yao et al., 2024) and Graph-of-127 Thought (Besta et al., 2024), aim to improve the search process within the *answer space*, rather 128 than the prompt space, and remain unsupervised. These "X-of-thought" approaches still rely on a 129 "one-prompt-for-all" strategy, where the model autonomously devises a step template (algorithm) 130 for each task. Once the template is established, these approaches help navigate the answer space 131 more effectively. For instance, Graph-of-Thought encourages the model to frequently revisit previously generated steps, while Tree-of-Thought allows the model to generate multiple possible next 132 steps before selecting the most promising one. However, the step template itself (algorithm) is still 133 generated by the model and can be poorly suited to the problem (Figure 1.a), especially when task-134 specific supervision (guidance) is lacking. 135

136 In this work, we thoroughly investigate the distinction between prompt space and answer space in 137 the CoT process. Building on insights from previous theoretical analyses of CoT (Li et al., 2024), 138 we explore why "supervision" is necessary and how it can be provided to guide the model in finding the optimal steps. We conduct extensive experiments on structured reasoning tasks, demonstrating 139 that task-specific "supervised" CoT is crucial for achieving optimal solutions and highlighting the 140 performance gap when supervision is used versus when it is not. Our work is the first of its kind 141 to focus on prompt space exploration and offers valuable insights into *understanding* and *designing* 142 effective prompt techniques for reasoning tasks. 143

144

119

120 121

145 146

147

148

149

2 DEMYSTIFYING COT: A STRAIGHTFORWARD UNDERSTANDING

In this section, we summarize key findings from previous theoretical analyses (Li et al., 2024; Zhang et al., 2024; Feng et al., 2024) of CoT prompting, presenting them in a unified and accessible manner. The conclusions drawn here will serve as a foundation for our subsequent analysis of supervised CoT.

150 151 152

153

2.1 LIMITATIONS OF TRANSFORMER ARCHITECTURE

154 Transformers, unlike recurrent networks, are not designed to perform reasoning over an arbitrary 155 number of sequential steps (depth) internally. Specifically, in a Transformer model, the hidden state 156 \mathbf{h}_{t-1} at time step t-1 is not reused when calculating \mathbf{h}_{t} (Figure 2.b), as it would be in recurrent 157 networks like RNN (Figure 2.a). Instead, the hidden state \mathbf{h} is passed forward only through the 158 layers of the Transformer (Dehghani et al., 2018) (Figure 1.c), not through time, which means 159 that the number of sequential steps is fixed and limited for any given Transformer architecture (Li et al., 2024; Zhang et al., 2024; Elbayad et al., 2019). In contrast, Recurrent Neural Networks 160 (RNNs) (Grossberg, 2013) allow the hidden state h to be passed through time steps via recurrent 161 connections (Figure 2.a), enabling sequential computation over h through an arbitrary number of input tokens. This capability allows RNNs to perform deeper reasoning over h, which is essential
 for solving complex tasks (Zhang et al., 2024).

The hidden state h plays a crucial role in reasoning, as it stores both reasoning memory and in-165 termediate reasoning results (Zhang et al., 2024). The ability to sequentially compute and update 166 h over time allows a model to build reasoning depth, which is necessary for addressing complex 167 problems. This depth advantage provided by recurrent connections cannot be replicated by autore-168 gressive models. Autoregressive models, instead of passing the hidden state \mathbf{h}_{t} forward, pass the 169 generated token y_+ . However, y cannot replace the role of h for the following reasons: y is a discrete 170 value extracted from h and only contains partial information (Figure 2.b), making it insufficient for 171 continued reasoning in many tasks. y exists outside the latent space where h operates (Figure 2.b), 172 meaning it cannot be used for computation in the same way that h can (Zhang et al., 2024).

173 174

2.2 NATURE OF REASONING

175 176

Reasoning inherently requires *sequential* depth. For tasks with input of length n, reasoning is typically performed step by step to arrive at the final result. Examples include counting (incrementing a counter iteratively), playing chess (updating the board state iteratively), and searching (marking visited nodes iteratively). To solve a given task, there is a theoretical lower bound on the required depth of computation (Sanford et al., 2024). Since models like Transformers can only perform a constant number of sequential reasoning steps over the hidden state h, they are unable to solve reasoning tasks where the depth requirement increases with the length of the input.

Consider chess as an example. For a sequence of chess moves, $\mathbf{x}_n = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$, to validate the n-th move, the n-th board state \mathbf{h}_n must be calculated. This requires n *sequential* computations, as the n-th board state depends not only on the sequence of moves \mathbf{x} but also on the previous board state \mathbf{h}_{n-1} . While a neural network could *memorize* the mapping from \mathbf{x}_n to the correct h (Arpit et al., 2017), bypassing the need for sequential computation, memorization is much more resource-intensive than reasoning. This is because memorization would require storing all possible *permutations* of \mathbf{x}_n and their corresponding resulted board states, an exponential challenge that eventually demands infinite memory to store instances of arbitrary length.

Thus, in the example of simulating a chess game, the model's internal representation h, which encodes the board state, must be sequentially computed n times to simulate the game. Transformers, which lack the infinite precision needed for memorization, cannot perform such tasks, as their hidden states h are computed a fixed number of times, regardless of the input length.

195 196 197

2.3 CoT + Autoregressive = Recurrent

199 As previous studies have shown Li et al. (2024); Zhang et al. (2024); Feng et al. (2024), Chain of Thought (CoT) effectively bridges the gap between autoregressive Liang et al. (2022); Liu et al. 200 (2022) models and recurrent structures Zhang et al. (2024) within large language models (LLMs). 201 Instead of merely outputting tokens to answer questions, CoT also generates intermediate steps 202 which are not part of the answers. These intermediate steps, represented as a sequence of natural 203 language tokens (o_1, o_2, \ldots, o_k) , act as a discretization of the latent information \mathbf{h}_n (Figure 2.c). 204 Given that natural language is a powerful medium for encoding nearly any type of information, h 205 is effectively transformed into a token sequence \mathbf{o} , which is then converted back into a vector \mathbf{h} 206 via the embedding layer. In this way, computational information is preserved through a process of 207 discretization followed by vectorization, represented as: $\mathbf{h}_t \xrightarrow{\text{discritization}} (\mathbf{o}_1, \mathbf{o}_2, \cdots, \mathbf{o}_k) \xrightarrow{\text{vectorization}}$ 208 \mathbf{h}_{t+1} (Figure 2.c). This approach, effectively achieve the same effect as $\mathbf{h}_t \Rightarrow \mathbf{h}_{t+1}$ in the RNN-like 209 recurrent network, allowing h to be recurrently updated by the network. 210

In the earlier chess example, the LLM generates intermediate reasoning steps as natural language strings during the CoT process. Specifically, it produces a sequence of tokens (e.g., in English) to describe the board state h_k after the first k moves, detailing the positions of pieces such as the bishop and the king. In the subsequent computation, the LLM reads this board description up to move k and uses it to calculate the k+1-th board state, thereby avoiding the need to re-compute the reasoning from scratch—something Transformers cannot do internally due to their non-recurrent architecture. 216 In conclusion, LLMs with CoT effectively extend the reasoning process from the model's internal 217 latent space \mathbb{H} to a natural language-based token space \mathbb{O} . Thanks to the powerful encoding ability 218 of natural language, intermediate reasoning steps are encoded and stored in text form, which the 219 model can reuse in subsequent computations. This approach significantly increases the model's rea-220 soning depth to T(n), where T(n) is the number of CoT steps performed. Under ideal theoretical conditions—such as infinite CoT steps and perfect information conversion between latent and text 221 space—LLMs with CoT can achieve Turing completeness, theoretically solving any problem, in-222 cluding those beyond symbolic tasks (e.g. recognizing regular languages). This theoretical analysis 223 provides strong guidance for designing effective "supervised" CoT approaches, which we introduce 224 in subsequent sections. 225

- 226
- 227 228 229

3 COT SEARCH SPACE = PROMPT SPACE + ANSWER SPACE

While theory suggests CoT-augmented LLMs can solve any problem Li et al. (2024), finding solutions in practice is much harder. CoT is limited by a finite number of steps, and the conversion from latent states h to token sequences o is imperfect. Consequently, only partial information is extracted at each step, making it crucial to identify the right data to continue the correct computation. We decompose the CoT reasoning into two components: template search within the prompt space and answer search within the answer space. We show how effective navigation of the prompt space can simplify answer space complexity and reveal limitations of unsupervised "X-of-thought" methods.

237 238

3.1 PROMPT SPACE

239 240

The latent vector **h** contains rich intermediate information when processing a task, including counters, sums, flags for binary indicators, and more. When LLMs are prompted to "think step by step" along with the task instance, they generate a *step template*, specifying which information from **h** to extract and discretize into tokens $(o_1, o_2, ..., o_k)$. Ideally, as $k \rightarrow \infty$ —meaning the length of the CoT is arbitrarily long—all vectorized information in **h** can be fully textualized, achieving *true* recurrence through autoregression. However, with limited k, only partial information is discretized.

If we define the amount of information stored in h as m 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 m-bit space. Thus, the total number of potential step templates is $C(m, s) = \frac{m!}{s!(m-s)!}$, which *estimates* the number of ways information can be extracted via CoT at each step. Each template defines an extraction of unique s bits of information.

For example, in the chess simulation case, h encodes details such as the <current board layout>, <the next player>, <board status>, <number of pieces taken by each player> and so on. When given the instruction to "think step by step", the model decides which information to extract based on the *step template* it generates. Extracting the wrong information might hinder reasoning in subsequent steps as *recurrence* can not be effectively performed on the needed information.

258 The prompt search complexity C(m, s) depends on both m, the total information in h, and s, the 259 amount of information each CoT step can extract. If a model is sufficiently trained, the total amount 260 of encoded information in h is proportional to the dimension size of h (Allen-Zhu & Li, 2023), d, 261 denoted by $m \propto d$. In this context, m represents the size of the search space, while s correlates with the length of CoT tokens o, as longer CoT steps tend to extract more information from h. Thus, 262 s serves as the search step size. In practice, step template search is not entirely random. Models 263 often find relevant templates using heuristics, which significantly reduces the search complexity of 264 C(m, s). However, identifying the optimal template remains challenging, and using an incorrect 265 template can severely degrade performance, as demonstrated in our experiments. 266

In conclusion, the *step (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, determining what information is needed at each step and how it should be used to compute the next state.

3.2 ANSWER SPACE

271

289

291 292 293

294

295

296

297

298

Once the model "decides" on the steps to follow during CoT, it performs reasoning accordingly. With a specific step (prompt) template p_i chosen from the prompt space \mathcal{P} , CoT iteratively executes $\mathbf{h}_t \stackrel{p_i}{\Longrightarrow} (\mathbf{o}_1^{(i)}, \mathbf{o}_2^{(i)}, \dots, \mathbf{o}_k^{(i)}) \Rightarrow \mathbf{h}_{t+1}$ to update **h** and calculate the next state, continuing this process until reaching the final state (solution). The complexity of finding solutions in the answer space depends on both the choice of \mathbf{p}_i and the nature of the task itself.

Each task embeds a different level of complexity in its answer space. For instance, in the chess simulation task of <finding a set of actions leading to game end>, the answer space $S = (s_1, s_2, ..., s_{\infty})$ contains all possible combinations of action sequences s. The solution set $C\mathcal{R} \subset S$ includes all valid action sequences that lead to the end of the game, being a subset of the entire answer space S. Solving the problem requires identifying one single correct action sequence $s_{correct} = (y_1, y_2, ..., y_T) \in C\mathcal{R}$.

If a fixed step (prompt) template for this task, such as $p_0 = \langle \text{extract current board} \rangle$ configuration at each step>, is used, the CoT process iteratively extracts the current board description and use it for calculating next board state in h to identify the valid next move y_i , eventually forming the correct answer $s_{\text{correct}} = (y_1, y_2, \dots, y_T)$. The complexity of navigating the answer space can be roughly measured by:

$$\frac{\operatorname{len}(\mathcal{CR})}{\operatorname{len}(\mathcal{S})} | p \tag{1}$$

This ratio measures the proportion of the solution space $C\mathcal{R}$ relative to the entire answer space S, given a specific template p. If the chosen template p extracts irrelevant information—such as determining which player is next at each step—the ratio simplifies to $\frac{\text{len}(C\mathcal{R})}{\text{len}(S)}$. In this case, each y_i would be generated randomly, as h can not be computed iteratively over useful information needed for extracting correct $y_t exttti$, making the correct answer only discoverable by *chance*.

Correctly identifying the step template p is crucial for reducing the complexity of $\frac{\text{len}(CR)}{\text{len}(S)} \mid p$, 299 300 as p dictates what information is recurrently overlayed in the process $\mathbf{h}_t \Rightarrow \mathbf{h}_{t+1}$ and in turn 301 what can be calculated, essentially acting as the "algorithm" for solving tasks in the CoT pro-302 cess. In the chess example, the optimal template would be <extract current board configuration at each step>, allowing the model to reason over the board state itera-303 tively, i.e., $\mathbf{h}_t \xrightarrow{\text{board state}} \mathbf{h}_{t+1}$. With the correct board state computed recurrently, the valid 304 305 next move y_{t} can be effortlessly derived from h_{t} . However, using a less relevant template, such as 306 <extract the number of pieces on the board at each step>, would expand the search space nearly to $\frac{\text{len}(CR)}{\text{len}(S)}$, as the number of pieces doesn't provide useful information for 307 308 determining the next valid move. Consequently, the model would have to recalculate the board state at each step from previously generated moves y_1 , which requires O(n) depth–Transformers, lim-310 ited by constant depth, cannot handle. As a result, the next action yt+1 would not benefit from the 311 CoT process.

312

313 314

315

3.3 COT AS AN UNSUPERVISED TASK SOLVER

316 CoT operates in an unsupervised manner for any given task, relying on a single universal prompt, 317 Think Step by Step, and leaving it to the model to generate its own step template $p \in \mathcal{P}$ for 318 extracting information at each step. Since humans do not supervise step completion, the generation 319 of steps—i.e., determining which information to extract from h and compute recurrently—comes 320 primarily from the model's heuristics. For example, in counting tasks, LLMs use learned heuristics to extract a *Counter* value from h and perform recurrent updates. However, these unsupervised, 321 heuristic-driven templates are often unreliable, as the model lacks the knowledge to identify key 322 components for computation, as demonstrated in previous work Valmeekam et al. (2022) and our 323 experiments.

324 3.4 COT VARIANTS AS UNSUPERVISED HELPERS FOR NAVIGATING ANSWER SPACE

326 In practice, the answer space S can be large and complex, and even with the optimal step (prompt) 327 template p, CoT can make errors. Various CoT variants, such as Tree-of-Thought (ToT) and Graphof-Thought (GoT), have been proposed to mitigate these mistakes in solution searching. While 328 these "X-of-thought" approaches don't dictate which specific information to extract at each step like p does, they improve solution finding by exploring multiple paths and self-verifying. For 330 instance, ToT explores multiple instances in the answer space simultaneously under some given 331 template p, unlike the single-path exploration of CoT. Specifically, information extracted from the 332 current hidden state ht using p is used to generate q possible answers for the next step, denoted 333 as $(y_{t+1}^{(1)}, y_{t+1}^{(2)}, \dots, y_{t+1}^{(q)})$. Each answer leads to a different next state h_{t+1} . In the example of 334 <finding a set of actions leading to game end>, the board state at step t is ex-335 tracted into descriptions using the correct template p and to form h_{t+1} , and instead of producing a 336 single next move y_{t+1} from h, multiple actions are derived. Each derived action along with previous 337 actions forms a unique path that leads to a potential solution in S. Since some paths may fail (e.g., 338 leading to a non-ending game), exploring multiple paths simultaneously increases the efficiency of 339 searching the answer space. The visualization is shown in Figure 3.



Figure 3: ToT mechanism. h_t is transitioned into different h_{t+1} , to explore more in *answer space*. How state is transitioned is dictated by the step template of CoT, which goes beyond what ToT offers.

Similarly, GoT improves search accuracy by iteratively revisiting previously generated *partial answers*. However, none of these approaches are supervised, as the model is not informed of the correct step template p and generates it on its own, extracting information at each step accordingly. X-of-Thought still relies on a "one-prompt-for-all" approach and only aids in finding answers after $p \in \mathcal{P}$ is fixed. As we have shown, this can lead to poor outcomes, since p directly influences the complexity of the answer space, and X-of-Thought may be too late to correct errors in some cases.

4 EXPERIMENTS

340 341

342 343 344

345 346

347

348

349 350

351

352

353

354

355

356 357

358

365

366

367 368

369

In this section, we conduct experiments to demonstrate the importance of supervision in the CoT process. Specifically, we design scenarios where the correct step template is provided through supervision, and compare them to cases where incorrect steps are simulated by the model. Our results show significant performance degradation when the step templates are incorrectly derived, highlighting the need for human supervision to ensure reliable task performance with LLMs.

The objective of our experiments is not to evaluate the reasoning performance of different LLMs, but to emphasize the critical role that "supervision" plays in CoT. Comparing the abilities of various models is beyond the scope of this work.

4.1 EXPERIMENTS DESIGNS

370 Although we used chess simulation as an example of reasoning with CoT due to its resemblance to 371 real-life complex reasoning tasks, tasks involving chess boards and actions can be difficult to im-372 plement and evaluate. Instead, we follow previous work Zhang et al. (2024); Delétang et al. (2022) 373 by focusing on more fundamental reasoning tasks for LLMs. Specifically, we evaluate tasks at three 374 levels of computability: Regular (R), Context-Free (CF), and Context-Sensitive (CS), each corre-375 sponding to tasks solvable by different levels of computational power, from deterministic automata all the way to linear bounded automata (restricted Turing machines). These tasks involve operations 376 such as counting, sorting, and number addition-basic operations that are required by more com-377 plex algorithmic problems (like NP problems). Each task has a strong dependency on identifying

378			1	Tana		TIM	Сот	Сот	Сот
379	Level	Task	RNN	RNN	Transformer	w/o CoT	Unsupervised	CR Supervised	IN Supervised
380		Modular Arithmetic	1.00	1.00	0.24	0.22	0.96	1.00	0.44
381	R	Parity Check	1.00	1.00	0.52	0.58	0.94	1.00	0.42
001		Cycle Navigation	1.00	1.00	0.62	0.50	0.78	1.00	0.26
382		Stack Manipulation	0.56	1.00	0.58	0.00	0.92	0.96	0.00
383	CF	Reverse List	0.62	1.00	0.62	0.00	0.80	0.96	0.38
384		Modular Arithmetic	0.41	0.95	0.32	0.00	0.82	0.94	0.50
385		Odds First	0.51	1.00	0.53	0.00	0.80	0.92	0.00
000	CS	Addition	0.50	1.00	0.54	0.00	0.84	0.88	0.00
380	CS .	Multiplication	0.50	0.59	0.52	0.00	0.14	0.44	0.00
387		Sorting	0.28	0.71	0.92	0.00	0.36	0.90	0.00
388									

Table 1: Main results across three levels of reasoning tasks. For LLMs without CoT, intermediate steps are explicitly prohibited in the prompt. In CoT generation, "CR Supervised" refers to when we provide the correct supervision. "IN Supervised" refers to when seemingly correct but suboptimal step templates are provided, simulating scenarios where the model makes mistakes in navigating the prompt space and derives incorrect step templates. Bolded numbers indicate performance greater than or equal to 0.9, while red indicates low performance (below 0.2). Results for RNN, Tape-RNN and Transformer are trained expert model by previous research (Delétang et al., 2022), they are solely used for reference and not compared with LLMs as it follows slightly different experiment settings.

	R			CF			CS			
Model	MA	PC	CN	SM	RL	MA	OF	AD	MU	SO
Unsupervised CoT	0.96	0.94	0.78	0.92	0.80	0.82	0.80	0.84	0.14	0.36
Unsupervised ToT	0.92	0.90	0.92	0.36	0.88	0.78	0.82	0.94	0.18	0.66
Unsupervised GoT	1.00	0.98	0.90	0.72	0.92	0.88	0.82	0.92	0.20	0.80
Correctly supervised CoT	1.00	1.00	1.00	0.96	0.96	0.94	0.92	0.88	0.44	0.90

Table 2: Variant of CoT in performing each task. Each task is named using the first two letters in Table 1.

406 407

405

389

390

391

392

393

394

395

396

408 the correct step template, thus allowing us to clearly observe the impact of selection on step template on CoT performance. 409

410 All of these tasks require a level of computability beyond the capabilities of the Transformer's inter-411 nal architecture Delétang et al. (2022). Specifically, they demand a minimum computational depth 412 that scales linearly with input length, surpassing the constant depth inherent to Transformer models. 413 Thus, solving these tasks necessitates the use of CoT, and correctly identifying the information to 414 extract during CoT is crucial for resuming computation and building the necessary depth.

415 We use GPT-4-o classic, a version that eliminates the use of external tools (e.g., calculators or pro-416 grams) and functions solely based on the model itself. We test each task using instances sampled 417 according to previous work (Zhang et al., 2024). To ensure that factors such as long-context infor-418 mation retrieval and tokenization do not affect the results, we follow the setup from prior research 419 and conduct controlled experiments. Details of our experimental design, including length sampling, 420 task specifications, format adjustments, and prompt usage, are provided in the Appendix.

421 We extend the previous findings on expert models Delétang et al. (2022), which are specifically 422 trained for particular tasks, to our experiments with LLMs. Due to differences in experimental 423 settings, the results from expert models are presented for reference rather than direct comparison. 424 Unlike prior research, which reports the best performance out of N trials Delétang et al. (2022); 425 Zhang et al. (2024) for each task instance, we report the average one-trail performance across all 426 tested instances. Our focus is on practical usability beyond the theoretical upper-bound computability analysis in previous work. The final results are shown in Table 1. 427

428

430

- 429 4.2 MAIN RESULT
- Recurrence is key for reasoning. As demonstrated in both expert models (RNN, Tape-RNN, and 431 Transformers) and LLMs, recurrence is the determining factor for *solving* tasks in each category.



Figure 4: Average success rate in deriving correct step template in each level of tasks.

Specifically, expert models like RNN and Tape-RNN show the ability to solve tasks across various categories with over 90% accuracy, depending on their memory architecture. Transformers, however, are limited by their shallow depth of reasoning, as shown earlier, and fail to solve any tasks. Similarly, LLMs without CoT, relying solely on internal Transformer reasoning, achieved 0% performance on most tasks, with low performance on others likely due to guessing. When CoT augments LLMs with recurrent computational power, accuracy improves significantly. These comparisons highlight the critical role of *recurrence* in a model's computability, reinforcing the analysis we previously discussed.

Role of Step Template in Reasoning Performance: Supervision Is Essential. We provide human supervision for all tasks, and we observed that, due to the relatively simple nature of the tasks, the model makes mistakes in finding the optimal step template less frequently. As a result, it is difficult to clearly observe the performance gap between optimal and non-optimal step templates. To address this, we introduce two types of supervision for each task: Correct Supervision (CR Supervised), where the model is guided with optimal steps to demonstrate the best possible performance, and Incorrect Supervision (IN Supervised), which simulates scenarios where the model derives incorrect steps to show how performance can degrade. We present examples of these supervised scenarios for each task in Table 3.

Task	Correct Supervision Example	Incorrect Supervision Example
Modular Arithmetic	Write down partial sums after each step	Write down paired sums of each two values
		at each step
Parity Check	Write down "even" or "odd" counter after	Write down whether the word is a target
	each word in each step	word at each step
Cycle Navigation	Write down which state you are in at each	Write down the total number of "forward"
	step	at each step
Stack Manipulation	Write down the resulting stack at each step	Write down the number of operations per-
		formed up to that step
Reverse List	Write down the partially reversed list after	Write down the value to be added to the
	each step from the back	reversed list and the remaining original list
Modular Arithmetic	Write down the formula with reduced val-	Write down the result of each performed
	ues in the performed operations at each	operation at each step
	step	

Table 3: Examples of correct and incorrect steps for performed reasoning tasks.

From Table 1, we observe that providing supervision yields noticeable improvements over the unsupervised "step-by-step" approach. Specifically, errors caused by the model's own derived step
templates are eliminated with correct supervision, resulting in better performance scores. In contrast, when the step template is intentionally set up incorrectly, we observe a *significant* performance
degradation, with some tasks performing as poorly as they would without using CoT.

To explain this further, when a step template is incorrectly specified (e.g., outputting the sum up to the current step for a task that requires counting appearances), the useful counter information c in h_{t} is not extracted. As a result, c is not carried forward into the next state h_{t+1} , leading to a failure in resuming the necessary calculations. While the wrongly specified information (e.g., the partial sum) is recurrently calculated, it does not lead to the correct final answer for the task. CoT Variants are Useful in Navigating Answer Space. We compare the results of different CoT variants for the same tasks. As shown in Table 2, both ToT and GoT improve performance over naive CoT. However, this improvement is due to correcting "incorrect calculations" during computation, not from improvements in step-template selection. ToT provides little benefit, as the tasks typically have only one path to the solution. In contrast, GoT shows greater accuracy gains, thanks to its self-revisiting mechanism,

492 **Prompt Space Analysis.** We further analyzed the model's performance in navigating the prompt 493 space, i.e., finding the correct (optimal) step template for each task. As shown in Figure 4, all 494 tasks involve relatively simple calculations, and the model exhibits a high average success rate in 495 identifying the correct template. Specifically, the success rate for R-type tasks exceeds 90%. As 496 task complexity increases, we observe a slight decline, with CS tasks showing an 84% success rate in extracting the correct information during CoT. We further include case studies showcasing how 497 "sub-optimal" steps are derived from unsupervised CoT process, which are shown in Figure 1 and 498 Appendix Figure 5, 6 and 7. 499

Lastly, we showcase how incorrect navigation in the prompt space leads to uncorrectable results.
As shown in Appendix Figures 8 and 9, the incorrect step template results in incorrect information
extraction, leading to a wrongly computed next state and ultimately increasing the difficulty of
searching the answer space.

504 505

5 SUPERVISED COT: USERS' PERSPECTIVE

506 507 508

5.1 How to Supervise?

509 As we've demonstrated, providing correct supervision is crucial for helping the model achieve ac-510 curate results. A natural question arises: how can effective supervision be derived? The key to good 511 supervision lies in understanding CoT's underlying mechanism, which essentially involves relaying 512 information through the text space. For tasks requiring multiple steps, users need to identify what 513 each step is and what key information should be extracted at each step. While this might seem 514 straightforward in the basic reasoning tasks used in our experiments, it becomes more complex for 515 challenging tasks, where correctly identifying the information requires careful task analysis. There-516 fore, human knowledge is critical for enhancing the model's computational abilities and can directly 517 influence task success. However, this supervision adds substantial workload, as each task demands a unique understanding of its computational structure. 518

519 520

521

5.2 WHEN TO SUPERVISE?

As we've observed, using an incorrect step template—whether model-derived or humaninjected—can result in significant performance degradation. Based on this, it's important to avoid
providing supervision unless you are reasonably confident that the steps will not hinder the reasoning
process. In cases of uncertainty, it may be better to rely on the model's own heuristics.

526 527

6 CONCLUSIONS

528 529

Our work offers a unique perspective on the mechanics of Chain of Thought (CoT) prompting and 530 its role in enhancing model reasoning. Through theoretical analysis and practical insights, we show 531 how CoT transforms latent information into text space, enabling iterative and resumable reasoning 532 steps that expand a model's computational depth. We further connect the model's problem-solving 533 capabilities with the complexities of finding solutions. Our analysis of prompt space and answer 534 space underscores the importance of identifying the correct step template to simplify navigation—an often overlooked aspect in prompt-related research. The success of CoT hinges not only on gener-536 ating steps but on extracting the right information at each stage. Our experiments demonstrate that 537 incorrect step templates can severely impact reasoning, reinforcing the importance of supervision. Even small errors in template selection can lead to significant failures. Our findings combine the-538 oretical analysis and experimental evidence, offering valuable insights into CoT's limitations and potential for improving reasoning tasks in large language models.

540 REFERENCES

547

556

558

562

563

564

565

566

567

571

572

573

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International conference on machine learning*, pp. 233–242. PMLR, 2017.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gian inazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of
 thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
 - George S Boolos, John P Burgess, and Richard C Jeffrey. *Computability and logic*. Cambridge university press, 2002.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. ACM Transactions on Intelligent Systems and Technology, 15(3):1–45, 2024.
 - Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.
 - Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. Neural networks and the chomsky hierarchy. *arXiv preprint arXiv:2207.02098*, 2022.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. Depth-adaptive transformer. *arXiv* preprint arXiv:1910.10073, 2019.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing
 the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- 577 578 Stephen Grossberg. Recurrent neural networks. *Scholarpedia*, 8(2):1888, 2013.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
 language models are zero-shot reasoners. *Advances in neural information processing systems*,
 35:22199–22213, 2022.
- ⁵⁸²
 ⁵⁸³ Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024.
- Jian Liang, Chenfei Wu, Xiaowei Hu, Zhe Gan, Jianfeng Wang, Lijuan Wang, Zicheng Liu, Yuejian
 Fang, and Nan Duan. Nuwa-infinity: Autoregressive over autoregressive generation for infinite
 visual synthesis. *Advances in Neural Information Processing Systems*, 35:15420–15432, 2022.
- Puyuan Liu, Xiang Zhang, and Lili Mou. A character-level length-control algorithm for nonautoregressive sentence summarization. Advances in Neural Information Processing Systems, 35:29101–29112, 2022.
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pp. 2847–2854. PMLR, 2017.

- 594 Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Transformers, parallel computation, and loga-595 rithmic depth. arXiv preprint arXiv:2402.09268, 2024. 596
- Aaditya K Singh and DJ Strouse. Tokenization counts: the impact of tokenization on arithmetic in 597 frontier llms. arXiv preprint arXiv:2402.14903, 2024. 598
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, 600 Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. Nature medicine, 29(8):1930-1940, 2023.
- Luke Thorburn and Ariel Kruger. Optimizing language models for argumentative reasoning. In 603 ArgML@ COMMA, pp. 27-44, 2022. 604
- 605 Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Large lan-606 guage models still can't plan (a benchmark for llms on planning and reasoning about change). In 607 NeurIPS 2022 Foundation Models for Decision Making Workshop, 2022.
- 608 Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. On the 609 planning abilities of large language models-a critical investigation. Advances in Neural Informa-610 tion Processing Systems, 36:75993-76005, 2023. 611
- 612 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
- 613 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-614 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language 615 models. arXiv preprint arXiv:2206.07682, 2022. 616
- 617 Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of large language models of code. In Proceedings of the 6th ACM SIGPLAN International Sympo-618 sium on Machine Programming, pp. 1–10, 2022. 619
- 620 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik 621 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Ad-622 vances in Neural Information Processing Systems, 36, 2024.
- Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1, 624 grade-school math and the hidden reasoning process. arXiv preprint arXiv:2407.20311, 2024. 625
- 626 Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. Don't trust chatgpt when 627 your question is not in english: A study of multilingual abilities and types of llms. arXiv preprint 628 arXiv:2305.16339, 2023.
- 629 Xiang Zhang, Muhammad Abdul-Mageed, and Laks V. S. Lakshmanan. Autoregressive + chain 630 of thought = recurrent: Recurrence's role in language models' computability and a revisit of 631 recurrent transformer, 2024. URL https://arxiv.org/abs/2409.09239. 632
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, 633 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv 634 preprint arXiv:2303.18223, 2023. 635
- 636 637

638

623

601

602

APPENDIX А

639 A.1 EXPERIMENTAL DESIGN 640

641 Our experimental setup carefully addresses potential pitfalls that could influence the model's per-642 formance, specifically focusing on tokenization and context length. Tokenization issues can sig-643 nificantly affect how models handle specific tasks, often leading to failures not tied to the model's 644 reasoning ability. To counter this, we reformatted task instances to eliminate tokenization biases. 645 Moreover, LLMs often struggle with retrieving information from long contexts, leading to hallucinations or forgotten data during extended reasoning processes. This tends to degrade accuracy, 646 as models fail to maintain accurate references to the initial task elements throughout longer se-647 quences. While these challenges are important in real-world applications of LLMs, they are outside

648	the scope of our investigation, which prioritizes analyzing the effect of using different step template.
649	To maintain a controlled environment, we restrict task lengths from 20 to 35 elements. This thresh-
650	old was determined from preliminary analysis, where longer task sequences often introduced issues
651	not related to reasoning but to the model's internal optimization process. When the task sequence
652	exceeds 35 steps, models can divide output over multiple contexts, which distorts accurate informa-
653	tion retrieval. By maintaining a manageable length, we isolate and evaluate the differences between
654	reasoning with and without CoT, avoiding disruptions caused by excessive context length. For each
655	task, we generate 50 instances using a pre-written script and the results are examined by humans.

656 657 TASK DESIGN

Each task involves simple rule-based iterations, emphasizing memory access and iterative processes.
The challenge for the model lies in its ability to execute these tasks within the constraints of its architecture and memory systems. Below, we describe each task in detail, including sample inputs and outputs. For the Regular (R) class tasks, we include the following:

663 TASK DESIGN

662

669

670 671

672

673

674

675

676

677

678

679 680

681 682

683

684

685

686

687

688

689

690 691

693

694

696

697

699

- 665 666 For the Regular (R) class tasks, we include the following:
- 667 668 TASK DESIGN

For the Regular (R) class tasks, we include the following:

- 1. Modular Arithmetic: Given a sequence of n numbers and basic operations (+, -), compute the result modulo 5. For example, the input 4 + 2 3 should yield 3.
- Parity Check: Determine if the word "banana" appears an even number of times in a list containing the words "apple" and "banana." For example, the input ("banana", "apple", "banana") yields True.
- Cycle Navigation: Based on a sequence of actions ("forward," "backward," "stay"), determine the final position in a 5-state cycle starting from state 1. For example, ("forward", "stay", "backward") will return state 1.

For the Context-Free (CF) class tasks, we use the following:

- Stack Manipulation: Given a list of fruit names representing a stack and a sequence of stack operations, compute the final stack. For example, applying (pop "banana", push "orange") to ("apple", "banana", "grape") results in ("apple", "orange", "grape").
 - 2. Reverse List: Reverse a list of vegetable names. For example, ("carrot", "potato", "onion") becomes ("onion", "potato", "carrot").
 - 3. Modular Arithmetic (Complex): Compute the result of an arithmetic expression with n operations modulo 5. For instance, $((2 + 4) \times (3 1)) \mod 5$ yields 0.

For the Context-Sensitive (CS) class tasks, we evaluate the following:

- Odd First: Extract all items at odd positions from a list of animal names, followed by those at even positions. For example, ("dog", "cat", "elephant", "tiger") yields ("dog", "elephant", "cat", "tiger").
- 2. Addition: Given two large numbers with n digits, calculate their sum. For instance, the input 123456 + 987654 yields 1,111,110.
- 3. **Multiplication:** Multiply two large numbers with n digits. For instance, the input 345×567 yields 195,615.
- 4. Sorting: Sort a list of integers using the insertion sort algorithm. For example, the input (8, 3, 5, 1) would result in (1, 3, 5, 8).



756	
757	${\mathscr O}$ calculate the new list with values at odd positions followed by even
758	positions in the following list: ['tree', 'river', 'car', 'dog', 'sun', 'moon',
759	'apple', 'cat', 'dog', 'cat', 'city', 'truit', 'banana', 'river', 'dot', 'moon', 'meat' ifruit' 'apple' idot' icat' ibanana' itree' isun' icity' iapple'
760	'meat', 'dog', 'fruit', 'tree', 'banana', 'dot', 'cat', 'meat', 'river', 'moon',
761	'city', 'dog', 'fruit', 'apple']. Do this step by step.
762	
763	(a) Lot's brook down the task step by step
764	Let's break down the task step by step.
765	Step 1: Identify Odd and Even Positions
766	In Python (or most programming languages), list indices start at 0. Therefore:
767	 Odd positions (1-based) correspond to even indices (0-based).
768	
769	Even positions (1-based) correspond to odd indices (0-based).
770	Step 2: List Elements at Odd Positions (0-based)
770	These correspond to even indices $(0, 2, 4, \text{ atc})$
773	
774	Let's extract those elements:
775	Index 0: 'tree'
776	 Index 2: loar!
777	
778	Index 4: 'sun'
779	Index 6: 'apple'
780	Index 8: 'dog'
781	Index 10: 'citv'
782	 Index 12: /hanana!
783	
784	Index 14: 'dot'
785	Index 16: 'meat'
700	Index 18: 'apple'
788	 Index 20: 'cat'
789	Index 22: 'tree'
790	
791	• Index 24: "apple"
792	Figure 6: Model-derived template: <write at="" down="" each="" element="" next="" position=""></write>
793	which is sub-optimal and leads to incorrect results. The partial list is not computed recurrently
794	during the step, and the final list is only generated at the end. Note that the full output is not shown
795	due to its length.
796	
797	
798	
799	
800	
801	



864	
865	calculate -9 - 10 - 14 - 2 - 1 + 13 - 3 - 8 + 12 + 9 - 8 + 15 - 8 - 4 +
866	10 + 10 + 1 - 4 + 5 - 10 - 12 + 2 + 9 - 11 + 5 + 15 - 12 - 13 + 12 - 12 - 11 - 1 + 2 + 8 - 14 - 9 + 6 - 11 + 8 + 7 - 2 do it step by step each
867	step i print the sum of 2i th value and 2i-1 th value
868	
869	A the break down the calculation star by star as not your request. Each star will add two
870	consecutive values from the list and print the result
871	consecutive values from the list and print the result.
872	The list of numbers is:
873	-9, -10, -14, -2, -1, 13, -3, -8, 12, 9, -8, 15, -8, -4, 10, 10, 1, -4, 5, -10, -12, 2, 9, -10, -12, -10, -12, -10, -12, -10, -12, -10, -10, -10, -10, -10, -10, -10, -10
874	
075	Now let's add two consecutive numbers at a time:
070 877	1. Step 1: $-9 + (-10) = -19$
878	2. Step 2: $-14 + (-2) = -16$
879	3. Step 3: $-1 + 13 = 12$
880	4. Step 4: $-3 + (-8) = -11$
881	5. Step 5: $12 + 9 = 21$
882	$\frac{1}{2} = \frac{1}{2}$
883	5. Step 6: $-6 + 15 = 7$
884	7. Step 7: $-8 + (-4) = -12$
885	8. Step 8: $10 + 10 = 20$
886	9. Step 9: $1+(-4)=-3$
887	10. Step 10: $5 + (-10) = -5$
888	11. Step 11: $-12 + 2 = -10$
889	12. Step 12: $9 + (-11) = -2$
890	13. Step 13: $5 + 15 = 20$
091	14. Step $14: -12 + (-13) = -25$
893	$45 \text{oter} \ 45 \ 10 \ (-12) = 0$
894	15. Step 15: $12 + (-12) = 0$
895	16. Step 16: $-11 + (-1) = -12$
896	17. Step 17: $2 + 8 = 10$
897	
898	Figure 8: Case study of performing simple arithmetic under a wrong template (intentionally pro- vided through supervision). Since the sumulative sum is not derived iteratively, the intended
899	value cannot be computed recurrently through CoT, leading to incorrect results. The correct step
900	template for this task should be <write calculated="" down="" th="" the="" to<="" total="" up="" value=""></write>
901	each step>. Note that the full output is omitted due to its length.
902	
903	
904	
905	
907	
908	
909	
910	
911	
912	
913	
914	
915	
916	
917	



969 results (intentionally provided through supervision). Since the stack status is not iteratively updated 970 and passed to the next state, the results cannot be tracked effectively. Tracking the total number of 971 items in the stack is not useful for deriving the final stack. The correct step template for this task 970 should be <write down the current stack status at each step>. Note that the 971 full output is omitted due to its length.