
Dynamic Chain-of-Thought: Towards Adaptive Deep Reasoning

Libo Wang
UCSI University
Nicolaus Copernicus University
free.equality.anyone@gmail.com

Abstract

1 To reduce the cost and consumption of computing resources caused by compu-
2 tational redundancy and delayed reward assignment in long CoT, this research
3 proposes the dynamic chain-of-thought with adaptive reasoning time and steps. The
4 researcher used simulation experiment to simulate the integration of D-CoT through
5 Python 3.13 IDLE combined with a Python simulator based on GPTs. At the same
6 time, the researcher used DeepSeek R1 as a control group to test and compare
7 the performance of the D-CoT simulator in processing MIT OpenCourseWare’s
8 linear algebra exam questions. Experimental results show that D-CoT is better
9 than DeepSeek R1 based on long CoT in three indicators: reasoning time, CoT
10 length (reasoning steps) and token count, which achieves a significant reduction in
11 computing resource consumption. In addition, this research has potential value in
12 deep reasoning optimization that is used as a reference for future dynamic deep
13 reasoning frameworks.

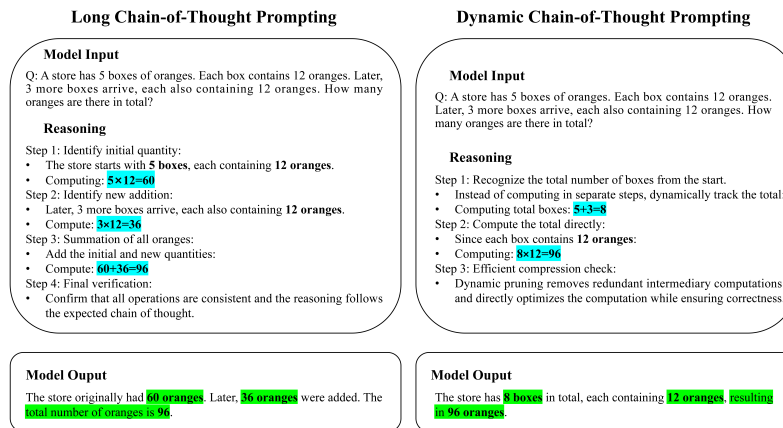


Figure 1: Comparison of the reasoning process of long CoT and dynamic CoT via prompting.

14 1 Introduction

15 As an emergent capability of large language models (LLMs) with huge parameter sizes in the
16 reasoning process, chain-of-thought (CoT) allocates additional computing resources in a way that
17 facilitates the gradual decomposition of tasks, which is particularly prominent in context learning [17].

18 It follows that chain-of-thought provides techniques for gradually unfolding intermediate reasoning
19 steps to enhance LLMs’ ability to handle complex problems by disassembling them into coherent
20 sub-steps [6]. However, affected by context, prompt word design, and model learning bias, CoT has
21 unfaithful interpretations, which leads to deviations between the reasoning process and the actual
22 decision-making mechanism [14]. In addition, while CoT enhances the reasoning capabilities of
23 LLMs, it also brings concerns about rising costs by prolonging the reasoning steps and improving
24 quality [8]. In response to the above-mentioned shortcomings of traditional CoT, more and more
25 cutting-edge LLMs apply long CoT technology to demonstrate excellent reasoning capabilities when
26 processing complex tasks [2, 16].

27 Long chain-of-thought (long CoT) aims to promote the model to have self-reflection and adaptive
28 refinement in multi-step complex scene reasoning through hierarchical reasoning and stepwise
29 verification to ensure accuracy and consistency [16]. For few-shot CoT, the accuracy of LLMs output
30 is linearly related to the number of reasoning steps, which means that the longer the number of steps,
31 the more accurate the response, while reducing the CoT length will significantly reduce the response
32 accuracy [8]. Jin et al. found that even if errors occur in the intermediate steps of long CoT during the
33 reasoning process, maintaining the necessary reasoning length will produce a high-accuracy response.

34 In the application of OpenAI’s o1 model, long CoT is combined with reinforcement fine-tuning
35 (RFT) to further optimize LLMs’ understanding and memory capabilities of multi-level reasoning
36 processes [7, 18]. In view of the shortcomings of LLMs in intermediate reasoning and adaptive
37 learning capabilities, DeepSeek-R1 introduces large-scale pure reinforcement learning training to
38 reduce reliance on supervised fine-tuning (SFT) [19]. During training, in order to avoid instability in
39 the initial cold start phase due to no longer relying on large amounts of labeled data, it constructs
40 long CoT data and uses specific collection and processing methods to guide deeper reflection, and
41 verification [2, 4].

42 However, long CoT requires a large number of intermediate steps in the reasoning practice process,
43 its inherent computational redundancy and delayed feedback problems significantly increase the
44 reasoning cost and consumption of computing resources, which is directly reflected in the exponential
45 growth of reasoning time and steps. In order to improve accuracy, a large number of lengthy reasoning
46 steps do not directly contribute to the final answer but only serve as an auxiliary process, resulting in
47 the accumulation of computational overhead [3]. These phenomena are often reflected in users’ actual
48 applications, especially LLMs with deep reasoning capabilities such as o3 min-high or DeepSeek R1.
49 For example, when users use DeepSeek R1 to perform difficult and complex tasks, the number of
50 reasoning steps increases significantly and the system response delay increases significantly.

51 The deep gap lies in the fact that long CoT essentially runs on a static reasoning framework, making
52 it difficult to flexibly adjust the number of reasoning steps according to the difficulty of different
53 problems. The lack of dynamic adaptation mechanism usually treats the reasoning process as a linear
54 expansion, which is reflected in the inability to adjust the length of the thinking chain according to
55 the complexity of different tasks and environmental feedback. For example, DeepSeek R1 faces the
56 problem of increased reasoning steps and longer system response delays when performing complex
57 and difficult tasks. The flaw makes it difficult for the system to adaptively allocate computing
58 resources, and the accumulation of inefficient reasoning causes nonlinear growth in reasoning time.
59 In addition, the insufficient coupling between long CoT generation and RL reward mechanism is one
60 of gaps about huge computational overhead. This technology is designed to enhance transparency
61 and explainability, but is not tightly integrated with RL goals, which results in the lack of an effective
62 credit allocation mechanism between reward signals during reasoning. Even though traditional RL
63 relies on immediate or deferred rewards to adjust strategies, the value of a step in a long CoT is
64 usually only determined when rewards are finally obtained.

65 **2 Proposed Framework & Algorithms**

66 In view of gaps, this research proposes dynamic chain-of-thought (D-CoT) to implement a state
67 compression mechanism with adaptive reasoning steps to reduce computational redundancy. Figure 1
68 shows the comparison of between long CoT and dynamic CoT under the same prompt.

69 **2.1 Dynamic Chain-of-Thought Framework**

70 Dynamic chain-of-thought (D-CoT) is an LLMs reasoning framework with adaptive reasoning
 71 capabilities that reduces the consumption of cost computing resources by real-time adjustment of
 72 chain length (reasoning steps) and reasoning time. Compared with the fixed and linear expansion of
 73 reasoning steps of traditional long CoT, D-CoT dynamically adjust the length of CoT in real time,
 74 select key steps after rating different tasks. Its specific internal structure is shown in Figure 2.

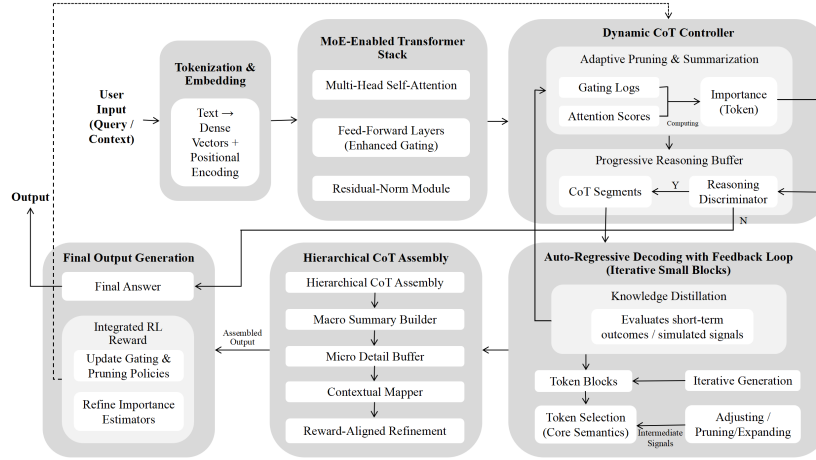


Figure 2: Dynamic Chain-of-Thought Framework

75 The D-CoT leverages the core principles of the hierarchical adaptive reinforcement learning to
 76 adjust the steps and information weights in the deep reasoning process to minimize computational
 77 redundancy and optimize the decision path. It introduces an importance-driven pruning strategy in
 78 the process of auto-regressive decoding, combines the partial reward estimator to instantly evaluate
 79 the effectiveness of the reasoning block, and decides whether to expand or delete the reasoning
 80 steps through adaptive thresholding. In addition, D-CoT constructs a multi-level reasoning structure
 81 through macro summary and micro detail buffer to ensure the optimal transmission of information
 82 flow at different reasoning scales. In contrast, it not only reduces the cumulative computational
 83 burden of long CoTs, but also improves the adaptability of reasoning, forming an efficient reasoning
 84 framework with feedback adjustment capabilities.

85 **2.1.1 Hierarchical Adaptive Reward Optimization**

86 The operating mechanism of the HARO (Hierarchical Adaptive Reward Optimization) algorithm
 87 is based on hierarchical reward allocation and adaptive reasoning adjustment. The algorithm uses
 88 the partial reward estimator to instantly evaluate the contribution of decisions at different levels of
 89 the deep reasoning process, and uses adaptive thresholding to dynamically correct the weight of the
 90 reasoning step (CoT length) to ensure the priority delivery of high-value information. In addition,
 91 this algorithm combines an importance-driven pruning strategy to instantly filter inefficient reasoning
 92 paths to reduce redundant computing overhead and improve overall reasoning efficiency.

93 **Token Importance Evaluation** Each reasoning step c_i is assigned an importance score:

$$I(c_i) = \alpha * A(c_i) + (1 - \alpha) * \text{GatingScore}(c_i)$$

94 where $A(c_i)$ is the dominance estimate derived from RL and GatingScore reflects the token-level
 95 contribution.

96 **Dynamic Adaptive Pruning Thresholding** It introduces a self-adjusting threshold τ_t based on
 97 historical success rates:

$$\tau_t = \gamma \cdot \tau_{t-1} + (1 - \gamma) \cdot \frac{1}{N} \sum_{j=1}^N 1[I(c_j)] > \tau_{t-1}$$

98 where τ_t represents the updated threshold at time step t ; γ is an attenuation factor used to control
 99 the retention of historical information; $1[I(c_j) > \tau_{t-1}]$ tracks whether past tokens have exceeded a
 100 previous threshold.

101 **Progressive Reasoning Buffer (Adaptive Selection)** Dynamic adjustments in CoT segments are
 102 stored in buffers:

$$C_t = C_{t-l} + \operatorname{argmax}_{c_i}(I(c_i) - \tau_t)$$

103 where steps below τ_t are discarded unless they contribute significantly to global coherence; C_t
 104 represents the CoT state at time t .

105 **Reward Optimization and Auto-Regressive Feedback** HARO uses reward gradients to iteratively
 106 optimize token selection by core semantics:

$$\nabla_{\theta} J = E [R_{\text{sem}}(C) + \lambda R_{\text{struct}}(C) \nabla_{\theta} \log \pi_{\theta}(C)]$$

107 where E represents expectation value; R represents reward function; θ represents model parameters;
 108 $R_{\text{sem}}(C)$ represents a semantic reward function; $R_{\text{struct}}(C)$ represents a structural function; γ is a
 109 weighting hyperparameter balancing semantic fidelity and structural efficiency; $\pi_{\theta}(C)$ represents the
 110 policy for selecting tokens; $\nabla_{\theta} \log \pi_{\theta}(C)$ is the optimization through policy gradients.

111 Notably, reward alignment and policy adjustment are inspired by PPO (Proximal Policy Optimization).
 112 PPO is committed to truncating the clipped objective function restriction policy and updating the
 113 amplitude policy update range to ensure training stability [12]. HARO further introduces adaptive
 114 thresholding, allowing reward signals to dynamically adapt according to the reasoning steps to
 115 improve the selectivity of the optimal decision trajectory. In addition, HARO draws on advantage
 116 estimation, dynamically filters high-value reasoning through hierarchical feedback mechanism,
 117 reduces low-reward expansion, and thereby reasons redundancy.

118 2.2 Detailed Framework Composition

119 The D-CoT framework consists of six key parts.

120 2.2.1 Tokenization & Embedding

121 As the first part of D-CoT, it is responsible for converting natural language input into a dense vector
 122 representation that can be processed by the model [13]. This process first decomposes the text into
 123 subword units through tokenization, maps it to a high-dimensional space through an embedding
 124 layer to capture semantic information and contextual dependencies, and combines it with positional
 125 encoding to provide sequence order information [15]. The following are the supported algorithms in
 126 the workflow:

Tokenization Process

$$\mathbf{T} = \text{Tokenizer}(\mathbf{Q})$$

Embedding with Positional Encoding

$$\mathbf{X} = E(\mathbf{T}) + P,$$

Seamless Transition to MoE Stack

$$\mathbf{X} \rightarrow \text{MoE-Enabled Transformer}$$

127 where converting user query \mathbf{Q} into a token sequence \mathbf{T} , this algorithm converts \mathbf{Q} into discrete
 128 labeled units; Mapping tokens into dense embeddings while integrating positional encoding; \mathbf{X}
 129 represents the final embedding representation that contains the vectorized representation and position
 130 encoding; $E(\mathbf{T})$ represents the embedding function that converts mark tokens into corresponding
 131 embedding representations; P represents position encoding.

132 When the processed dense vector representation enters the MoE-enabled transformer stack, the system
 133 selects the appropriate expert path through the gating network based on the semantic features of
 134 the input sequence that allows different parts of the information to flow through the most suitable
 135 expert layers (Pan et al., 2024). $\mathbf{X} \rightarrow \text{MoE-Enabled Transformer}$ means passing the final embedding
 136 representation \mathbf{X} to the MoE-enabled transformer stack.

137 **2.2.2 MoE-Enabled Transformer Stack**

138 The MoE-enabled transformer stack uses a mixture of experts (MoE) mechanism to enhance selective
 139 computing capabilities through multi-head self-attention to ensure efficient acquisition of key infor-
 140 mation [4]. The feed-forward layers combined with enhanced gating perform feature transformation
 141 based on dynamically selected experts to maximize reasoning efficiency [19]. The residual-norm mod-
 142 ule provides gradient stability and reduces signal attenuation to ensure smooth flow of information in
 143 deep structures. The researcher demonstrated its workflow algorithm as follow.

Expert Selection

$$\alpha_{t,e} = \text{Router}(u_t, e), E_{\text{active}} = \text{TopK}\{\alpha_{t,e} \mid e = 1, \dots, N\}$$

Expert Outputs

$$h'_t = \sum_{e \in E_{t,\text{active}}} \alpha_{t,e} \cdot f_e(u_t), u_t \in \mathfrak{X}$$

Importance Score

$$I_t = \gamma \left(\sum_{e \in E_{t,\text{active}}} \alpha_{t,e} \right) + (1 - \gamma) \beta_t$$

144 After expert assignment of input vectors through MoE, the refined semantic information is passed to
 145 the dynamic CoT controller for adaptive pruning and dynamic summarization

Adaptive Thresholding for Pruning and Summarization

$$\tau_{\text{dyn}}(r_t) = \tau_0 + \eta(r_t - \bar{r})$$

Hierarchical Decoding & Assembly

$$y_{t+1} = \text{Assemble}(B_T, \{\text{macro}, \text{micro}\})$$

146 where $\tau_{\text{dyn}}(r_t)$ is dynamic threshold; r_t is partial reward; τ_0 is a base threshold, \bar{r} is a running
 147 average reward, and η is a scaling factor; y_{t+1} represents the next generated token after decoding; B_T
 148 represents the final buffer of CoT segments after iterative refinement; The macro-level summaries
 149 compress global information; micro-level expansions retain fine-grained details.

150 **2.2.3 Dynamic CoT Controller**

151 As the core of deep reasoning, the dynamic CoT controller is responsible for optimizing the length
 152 and content of CoT to reduce redundant calculations and the accumulation of unnecessary steps.
 153 The key technology is to dynamically adjust the reasoning steps to adapt the reasoning process
 154 to different types of task requirements, thereby improving the adaptability of LLMs in complex
 155 decision-making scenarios. Its operating mechanism calculates the importance score of the token
 156 based on gating logs and attention scores, which determines whether the inference step should be
 157 retained or compressed. Afterwards, the reasoning fragments are stored in the progressive reasoning
 158 buffer to ensure that key information is efficiently used in subsequent steps. Reasoning discriminator
 159 determines whether to start CoT reasoning through knowledge certainty evaluation (Y/N). The system
 160 calculates the confidence $P_{\text{fact}}(x)$ via FAISS/BM25, $C_{\text{comp}}(x)$ is not greater than 3 that is based
 161 on syntactic structure and computational cost evaluation. If both are below the threshold, they are
 162 directly output, otherwise they enter CoT segments for reasoning. The relevant supported algorithms
 163 are as follows:

Adaptive Pruning & Summarization

$I_t = \text{Importance}(t)$ (combining gating + attention)
 if $I_t < \tau_{\text{dyn}}(r_t)$, prune token t , else, optionally summarize

Progressive Buffer Update

$$B_{t+1} = B_t \cup \{\text{Summarized } t \mid I_t \geq \tau_{\text{dyn}}(r_t)\}$$

Partial Reward

$$r_t = \text{RLFeedback}(t), \tau_{\text{dyn}}(r_t) = \tau_0 + \eta \cdot (r_t - \bar{r})$$

Adaptive Token Expansion and Pruning

$$z_{t+1} = \text{Adjust}(\text{Generate}(\text{SelectTokens}(B_t, \theta_t, r_t), B_t, \theta_t), r_{t+1})$$

Reasoning Discriminator

$$y_{\text{out}} = \begin{cases} \text{OutputAnswer}(x), & \text{if } P_{\text{fact}}(x) \geq 0.85 \text{ and } C_{\text{comp}}(x) \leq 3 \\ \text{CoT Segments}(x), & \text{otherwise} \end{cases}$$

Reward-Guided CoT Assembly

$$B_t = \text{AssembleCoT}(B_t, \{z_t\}_{t=-T}, r_{t+1} = \text{RewardUpdate}(r_t, \text{DecodeOut}(z_{t+1})))$$

164 where I_t is importance score of token t (gating + attention); τ_{dyn} represents dynamic threshold
165 based on reward r_t ; τ_0 is base pruning threshold; η is scaling factor for threshold adjustment; r_t is
166 RL feedback at step t / partial reward signal; r_{t+1} represents the updated reward at step $t+1$; \bar{r}
167 represents running average reward; B_t represents token buffer at step t / the final structured CoT
168 assembly after multiple iteration steps; B_{t+1} represents updated token buffer; t represents the current
169 reasoning step in the iterative decoding process; Summarized t is compressed token representation;
170 $\text{RLFeedback}(t)$ represents RL-based reward function for token t ; z_t represents tokens at time step
171 t / the generated reasoning tokens at step t ; θ_t represents hidden parameters of the decoder at t ;
172 $\text{SelectTokens}()$ represents token selection of core semantics; $\text{Generate}(\cdot)$ is based on generation of
173 new tokens; $\text{Adjust}(\cdot)$ expands or prunes tokens (adaptive pruning or expansion decision) based on
174 rewards and gating logs; T is The total number of CoT reasoning steps; $\text{RewardUpdate}()$ updates the
175 reward based on the previous reward and newly decoded token sequence; $\text{DecodeOut}(z_{t+1})$ is from
176 the iterative reasoning process at step $t+1$; $\text{AssembleCoT}()$ structures the final CoT by integrating
177 generated reasoning tokens and their associated reward values.

178 The processed CoT fragment is passed to auto-regressive decoding with feedback loop. It evaluates
179 the effectiveness of each fragment based on the partial reward estimator and optimizes the reasoning
180 strategy through dynamic adjustment, which reduces computational overhead while maintaining high
181 accuracy output.

2.2.4 Auto-Regressive Decoding with Feedback Loop (Iterative Small Blocks)

183 Auto-regressive decoding with feedback loop calculates the relative contribution of tokens in the
184 current reasoning process through the partial reward estimator. It filters qualified tags into iterative
185 generation and gradually builds a complete reasoning chain. The generation process uses small block
186 decoding to ensure that reasoning can optimize the structure under short-term feedback constraints
187 to improve adaptability and efficiency. Then, adjusting/pruning dynamically selects whether to
188 compress, retain or expand tags based on the partial reward signal and internal dynamic threshold
189 to reduce redundant calculations and ensure efficient reasoning. Some of the filtered tokens are fed
190 back to the progressive reasoning buffer to ensure consistent reasoning logic, and weight estimation
191 is dynamically updated. The mathematical expression of attention merging is:

Hierarchical Decoding and Expansion Rule

$$C_{t+1} = C_t \cup T_{t+1}, T_{t+1} = \text{Adjust}(\text{Decode}(T_t, \Theta_t, r_t), \text{Importance}(T_{t+1}), r_{t+1})$$

Reward-Guided Hierarchical CoT Refinement

$$F(t+1) = \text{Refine}(\text{AssembleSplit}(T_{t+1}), \text{RewardMap}(M_{t+1}), r_{t+1})$$

192 where C_t is CoT structure at time step t ; C_{t+1} represents update CoT after processing new tokens;
193 T_t is token block at time step t ; T_{t+1} is adjusted token block after pruning/expansion; θ_t is Local
194 model parameters, including gating logs and RL signals; r_t represents partial reward signal guid-
195 ing RL; r_{t+1} represents partial reward signal guiding hierarchical selection; $\text{Importance}(T_{t+1})$
196 represents function computing token importance; $\text{Adjust}()$ is the operator for pruning/expanding
197 token sequences; $\text{Decode}()$ represents function generating token sequences from parameters; F_{t+1} is

198 final hierarchical CoT representation at step $t + 1$; $\mathcal{M}_t + 1$ is macro-level summary tokens, $m_t + 1$ is
 199 micro-level detail tokens; AssembleSplit() represents operator splitting token blocks into macro/micro
 200 segments; RewardMap() represents function ranking macro segments based on reward signals; Re-
 201 fine() represents operator merging macro, micro, and reward-driven structures into final hierarchical
 202 CoT.

203 2.2.5 Hierarchical CoT Assembly

204 Hierarchical CoT assembly extracts high-level semantic information through macro summary builder
 205 and establishes a global reasoning structure to ensure contextual consistency. The fine-grained
 206 information is stored through micro detail buffer, which retains important reasoning details to avoid
 207 semantic loss. The stored information is further analyzed by the contextual mapper, which enables
 208 the reasoning process to dynamically adapt to contextual changes by integrating different levels of
 209 information. The reward-aligned refinement adjusts the reasoning weight through the RL mechanism
 210 and dynamically adjusts the contribution of key markers based on the learning reward signal. The
 211 algorithm is as follows:

Macro / Micro Separation

$$C_{\text{macro}}, C_{\text{micro}} = \text{Assemble}(C, \text{PartialRewards})$$

Reward-Aligned Refinement

$$C_{\text{final}} = \text{Refine}(C_{\text{macro}}, C_{\text{micro}}, \text{RewardMap})$$

212 where C still represents CoT; C_{macro} represents is a macro-level reasoning component; C_{micro} is
 213 micro-level reasoning component; Assemble(\cdot) is assembly function, it is responsible for dividing the
 214 reasoning steps and decomposing the complete CoT structure into C_{macro} and C_{micro} ; C_{final} represents
 215 final refined CoT; Refine(\cdot) represents refinement function.

216 2.2.6 Final Output Generation

217 After completing the hierarchical CoT assembly, the final output generation produces the final answer
 218 in response. At the same time, integrated RL reward re-evaluates reasoning weights and adjusts
 219 pruning and gating policies through importance estimators, and feeds the updated decision signals
 220 back to the dynamic CoT controller to optimize future reasoning processes. The supported algorithms
 221 are as follows:

Final Answer

$$\mathbf{A} = \text{OutputAnswer}(C_{\text{final}})$$

Compute Episode Reward

$$R_{\text{episode}} = \text{RewardFunction}(\mathbf{A}, \text{EnvironmentState})$$

Update Dynamic CoT Parameters (Loop)

$$\theta_{\text{CoT}} \leftarrow \theta_{\text{CoT}} + \eta \nabla_{\theta_{\text{CoT}}} R_{\text{episode}}$$

222 where \mathbf{A} is final answer; C_{final} represents the output result from hierarchical CoT assembly; R_{episode}
 223 represents episode reward; θ_{CoT} is a parameter of D-CoT; $\nabla_{\theta_{\text{CoT}}} R_{\text{episode}}$ is gradient of episode
 224 reward; η is the learning rate that controls the update amplitude of the D-CoT parameters.

225 3 Experiments

226 D-CoT is essentially a modular enhancement technology for deep reasoning. Its core mechanism
 227 includes adaptive step adjustment and reasoning pruning, but does not involve changes to the pre-
 228 training architecture of LLMs. It needs to be embedded into existing LLMs to achieve performance as
 229 it relies on its language understanding capabilities rather than independent execution [9]. Therefore,
 230 this research chose simulation experiment to test D-CoT’s adaptive step adjustment, computational
 231 resource allocation, reasoning pruning and reward alignment in the deep reasoning process. It allows

232 verifying the impact of different regulation strategies on reasoning performance in a controlled
233 environment , which avoids being limited by the immutability of the internal architecture [5, 10].

234 In addition, except for a few open sources such as DeepSeek R1 and MemGPT, mainstream ones such
235 as OpenAI o1 and o3 min-high are still closed-source LLMs, and the API interface cannot provide
236 internal control capabilities for the reasoning steps [11]. The researcher was unable to directly modify
237 the internal reasoning architecture of LLMs by integrating D-CoT for dynamic step adjustment and
238 reward alignment testing [1].

239 3.1 Experimental Setup

240 Considering that it is difficult to directly integrate into existing closed-source LLMs and a simulation
241 experiment was used to test its reasoning performance, the researcher simulated the process of
242 integrating D-CoT into current LLMs through custom GPTs with runnable code. By developing
243 code in Python 3.13 IDLE and uploading it to the Python simulator based on GPTs, the researcher
244 simulate the dynamic control of D-CoT’s deep reasoning steps, calculation mechanism adjustment
245 and reward alignment strategies, and simulate its time application scenarios in LLMs. The Python
246 simulator has the highly complex ability to execute code in customized GPTs, which earned it a
247 4.2-star rating, ensuring flexibility and controllability in testing. The researcher uploaded the D-CoT
248 code to a Python-based simulator as an experimental group, and selected DeepSeek R1 with deep
249 reasoning capabilities as a control group to compare the performance differences of reasoning.

250 3.2 Dataset

251 Based on the fact that solving linear algebra requires high structure and strong reasoning requirements,
252 it has become an effective test set to evaluate the reasoning ability of CoT. The researcher chose
253 the test questions of 18.06 Spring 2022 Problem Sets and Exams of 18.06 Linear Algebra of MIT
254 OpenCourseWare as experimental data. Linear algebra usually includes multi-step calculations and
255 logical derivation. Both Long CoT and D-CoT need to show the details in the reasoning process. In
256 addition, the test questions cover questions of different difficulty levels and can test the adaptability
257 and robustness of dynamic D-CoT under various reasoning challenges. And because this test question
258 is provided in text format, it is suitable for uploading to GPTs that emulate Python and ensures the
259 independent operation of D-CoT. This data has been placed in Appendix 1, and it is noteworthy that
260 the MIT OpenCourseWare license authorizes the test questions to be used publicly and experimentally
261 for non-commercial purposes.

262 3.3 Implementation

263 The researcher uploaded the D-CoT code developed by Python 3.13 IDLE to the Python simulator
264 based on GPTs to simulate its integration process in LLMs. At the same time, the researcher also set
265 up corresponding instructions in DeepSeek R1 to execute and process linear algebra test questions.
266 Afterwards, the researcher uploaded the 18.06 Spring 2022 Problem Sets and Exams one by one to
267 the simulator and DeepSeek R1 integrated with D-CoT and recorded the experimental results. In
268 order to avoid deviations in image semantics recognition between the experimental group and the
269 control group, all test questions were converted to machine-readable format and renumbered into 18
270 questions without changing the original meaning.

271 In view of the problem, the researcher sets three core evaluation indicators to evaluate the reasoning
272 cost and computing resource consumption of D-CoT. These metrics are reasoning time, CoT length
273 (reasoning steps) and token count to quantify the computational overhead of the D-CoT simulator and
274 DeepSeek R1 respectively. In addition, since this research is committed to reducing computational
275 redundancy and optimizing reasoning steps rather than verifying the accuracy of calculation results,
276 it does not include accuracy score as an evaluation indicator. Complete experimental records and
277 codes have been uploaded to supplementary material to ensure reproducibility.

278 3.4 Result & Discussion

279 After the above execution process, the experimental group and the control group respectively display
280 the data results except for reasoning time, CoT length (reasoning steps) and token count (Figure 3).

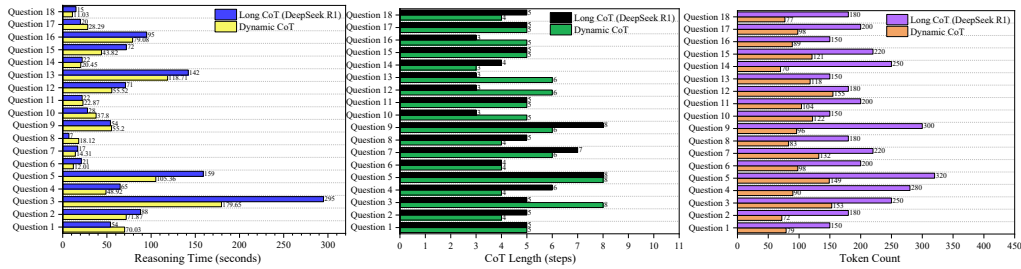


Figure 3: Comparison between D-CoT simulator and DeepSeek R1

281 According to the data results, in terms of reasoning time, the maximum reasoning time of DeepSeek
 282 R1 reaches 295 seconds, while D-CoT is significantly reduced to 179.65 seconds, showing its
 283 computing optimization capabilities in long reasoning processes. Judging from the distribution trend,
 284 the reasoning time of DeepSeek R1 is mostly concentrated in 50 to 150 seconds, but some questions
 285 exceed 200 seconds, showing the instability of its reasoning chain length; in comparison, D-CoT
 286 mostly maintains in the range of 40 to 120 seconds, and the fluctuation is small, indicating that it can
 287 effectively control the reasoning time. In terms of CoT length (reasoning steps), the number of steps
 288 of DeepSeek R1 is up to 8 steps, while D-CoT is controlled within 6 steps, and most questions only
 289 require 3 to 6 steps, which significantly shows that it can reduce redundant reasoning steps through a
 290 dynamic adjustment mechanism. As for the token count, the token usage of DeepSeek R1 is up to
 291 320 tokens, while that of D-CoT is significantly reduced to 180, and its overall distribution is mainly
 292 concentrated in the 70 to 180 range. Compared with DeepSeek R1, which exceeded 300 in some
 293 complex questions, D-CoT has demonstrated effective suppression of token growth and reduction of
 294 computing resource consumption.

295 4 Limitation & Future Research

296 In light of the fact that external researcher do not have the authority to access and modify the internal
 297 architecture of most current mainstream closed-source LLMs, D-CoT is difficult to directly integrate
 298 and uses simulation experiments to compare with DeepSeek R1 through a customized GPTs platform.
 299 However, the computing power, neural network size and parameters of GPTs are significantly different
 300 from DeepSeek R1, which causes the generalizability of the data results to be affected by additional
 301 factors. In addition, D-CoT relies on auto-regressive decoding and feedback mechanisms to adjust
 302 dynamic reasoning steps. This mechanism causes certain interference in fully reproducing the effects
 303 of pruning and reward alignment of LLMs under different contextual conditions in a simulation
 304 environment. The above limitations show that obtaining LLMs of open source architecture for
 305 internal integration testing in the future will help to more accurately evaluate the actual application
 306 performance of D-CoT.

307 5 Conclusion

308 In view of the fact that long CoT has a large number of intermediate steps in reasoning, which
 309 increases reasoning costs and consumption of computing resources due to inherent computational
 310 redundancy and delayed feedback, this research proposes dynamic chain-of-thought (D-CoT) to
 311 implement a dynamic deep reasoning mechanism of adaptive pruning, reward alignment, and step-
 312 wise control. It simulated the integration of D-CoT through simulation experiment in Python 3.13
 313 IDLE combined with the Python simulator based on GPTs, and used DeepSeek R1 based on long CoT
 314 as a control group to test its performance in solving the 18.06 linear algebra test questions of MIT
 315 OpenCourseWare. Experimental results show that D-CoT can effectively reduce computing resource
 316 consumption in reasoning time, CoT length and token count, thereby optimizing and reducing the
 317 computing cost and resource consumption of deep reasoning. Although D-CoT is limited by the
 318 current closed-source environment and cannot be directly integrated into LLMs for testing, it provides
 319 a new feasibility perspective for deep reasoning optimization of LLMs.

320 **Appendix 1**

321 The data for this research comes from the 18.06 Spring 2022 Problem Sets and Exams questions of
 322 MIT OpenCourseWare’s 18.06 Linear Algebra, which were uploaded to the D-CoT simulator and
 323 DeepSeek R1 to answer respectively. In order to avoid the image semantic recognition bias of matrix
 324 operations from interfering with the experimental results, the researcher converted all test questions
 325 into computer language input and renumbered them into 18 questions without changing the original
 326 meaning. This test may be used non-commercially for the experimental purposes of this study under
 327 the terms of the Creative Commons Attribution-Non Commercial sharealike (CCBY-NC-SA) license.

328 **Problem 1** (4+4+6 points):

329 The matrix A is given by $A = LUL^{-1}U^{-1}$

330 for

$$L = \begin{pmatrix} 1 & & & \\ -1 & 1 & & \\ 0 & 3 & 1 & \\ 1 & 0 & 0 & 1 \end{pmatrix}, U = \begin{pmatrix} 2 & 0 & 1 & 1 \\ & -1 & 0 & -1 \\ & & -2 & 1 \\ & & & 1 \end{pmatrix}.$$

331 (a) Write an expression for A^{-1} in terms of L, U, L^{-1} , and/or U^{-1} (but you don’t need to actually
 332 multiply or invert the terms!).

333 (b) What is the determinant of A?

334 (c) Solve $PAx = b$ for x , where P is the 44 permutation that swaps the 1st and 4th elements of a

335 vector, and $b = \begin{pmatrix} -5 \\ 4 \\ 11 \\ -3 \end{pmatrix}$. (You can get partial credit by just outlining a reasonable sequence of steps
 336 here that doesn’t involve a lot of unnecessary calculation.)

337 **Problem 2** (4+6 points):

338 (a) If a and x are vectors in \mathbb{R}^n , then $aa^T x$ can be computed using either left-to-right as $(aa^T)x$ or
 339 right-to-left as $a(a^T x)$, where the parentheses indicate the order of operations. Roughly count the
 340 number of arithmetic operations (additions and multiplications) in these two approaches: say whether
 341 each approach scales proportional to n, n^2, n^3 , etcetera.

342 (b) A is an $n \times n$ real matrix and x is an n -component real vector. Indicate which of the following
 343 must be equal to one another:

344 $\text{trace}(Axx^T), \text{trace}(xAx^T), \text{trace}(x^T Ax), x^T Ax,$
 345 $\text{trace}(x^T xA), xx^T A, \text{trace}(xx^T A), \text{determinant}(xx^T A).$

346 For the expressions that are equal, indicate how you would evaluate this quantity in a cost (in
 347 arithmetic operations) proportional to n^2 .

348 **Problem 3** (4+4+4+5 points):

349 You have a 4×3 matrix $A = (q_1 \quad 2q_2 \quad 3q_1 + 4q_2)$, where we have expressed the three columns of
 350 A in terms of the orthonormal vectors

$$q_1 = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}, \quad q_2 = \frac{1}{2} \begin{pmatrix} 1 \\ -1 \\ -1 \\ 1 \end{pmatrix}.$$

351 (a) What is the rank of A?

352 (b) Give a basis for $N(A)$.

353 (c) You are asked to calculate the projection matrix P onto $C(A)$. Your friend Harvey Ard suggests
 354 applying the formula $P = A(AT A)^{-1}AT$ he memorized in linear algebra. Explain why this won’t work
 355 here, and give an even simpler (correct) formula for P in terms of the quantities above. (You need not
 356 evaluate P numerically, just write a formula in terms of products of quantities defined above.)

357 (d) Find the closest vector to $x = \begin{pmatrix} 2 \\ 0 \\ 0 \\ 0 \end{pmatrix}$ in $N(A^T)$.

358 **Problem 4** (3+4+4+6 points):

359 The nullspace $N(A)$ of the real matrix A is spanned by the vector $v = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix}$

360 (a) Give as much true information as possible about the size (the number of rows and columns) of A .

361 (b) Give an eigenvector and eigenvalue of the matrix $B = (3I - A^T A)(3I + A^T A)^{-1}$.

362 (c) Aside from the eigenvalue identified in the previous part, all other eigenvalues λ of B must be
363 (circle/copy all that apply): purely real, purely imaginary, zero, negative real part, positive real part,
364 $j\lambda_j < 1, j\lambda_j > 1, j\lambda_j \leq 1, j\lambda_j \geq 1$.

365 (d) Give a good approximate formula for B^n

$$\begin{pmatrix} 0 \\ -1 \\ 0 \\ 8 \end{pmatrix}$$

366 for large n . (Give an explicit numerical vector, possibly including simple functions of n like 2^n or
367 n^3 ... no other abstract symbolic formulas.)

368 **Problem 5** (10 points):

369 Describe (give an explicit numerical result with as few unknowns as possible) all possible linear
370 combinations of the vectors

$$a_1 = \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix}, a_2 = \begin{pmatrix} 1 \\ 2 \\ 4 \end{pmatrix}, a_3 = \begin{pmatrix} 1 \\ -1 \\ 3 \end{pmatrix}, a_4 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

371 that give the vector $x = \begin{pmatrix} 4 \\ -1 \\ 5 \end{pmatrix}$

372 **Problem 6** (8+8 points):

373 Professor May Trix is trying to construct an 18.06 homework question in which $dx/dt = Ax$ has the
374 solution

$$x(t) = v_1 \cos(2t) + v_2 e^{-t} + v_3 \sin(2t)$$

375 for some **nonzero real** constant vectors v_1, v_2, v_3 , and some initial condition $x(0)$. Help May
376 construct A, v_1, v_2, v_3 , and $x(0)$:

377 (a) Write down a numerical formula for a possible real matrix A such that A is as small in size as
378 possible and where A contains no zero entries. Your formula can be left as a product of some matrices
379 and/or matrix inverses - you don't need to multiply them out or invert any matrices, but you should
380 give possible numeric values for all of the matrices in your formula. (You don't need to explicitly
381 check that your A has no zero entries as long as zero entries seem unlikely. e.g. the inverse of a
382 matrix with no special structure probably has no zero entries.)

383 (Note that there are many possible answers here, but they will all have certain things in common.)

384 (b) Using the numbers that you chose from the formula in your previous part, give possible corre-
385 sponding (numeric) values for $x(0), v_1, v_2$, and v_3 .

386 **References**

- 387 [1] Aitor Arrieta, Miriam Ugarte, Pablo Valle, José Antonio Parejo, and Sergio Segura. O3-mini vs
388 deepseek-r1: Which one is safer?, 2025.
- 389 [2] Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi
390 Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Do
391 not think that much for $2+3=?$ on the overthinking of o1-like llms, 2025.
- 392 [3] Damai Dai, Chengqi Deng, Chenggang Zhao, R.x. Xu, Huazuo Gao, Deli Chen, Jiashi Li,
393 Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y.k. Li, Panpan Huang, Fuli Luo, Chong
394 Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization
395 in mixture-of-experts language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar,
396 editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*
397 *(Volume 1: Long Papers)*, pages 1280–1297, Bangkok, Thailand, 2024. Association for
398 Computational Linguistics.
- 399 [4] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, and et al. Deepseek-v3 technical report, 2025.
- 400 [5] BRUCE EDMONDS and DAVID and HALES. Computational simulation as theoretical experi-
401 ment. *The Journal of Mathematical Sociology*, 29(3):209–232, 2005.
- 402 [6] Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards
403 revealing the mystery behind chain of thought: A theoretical perspective. *Advances in Neural*
404 *Information Processing Systems*, 36:70757–70798, 2023.
- 405 [7] Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie
406 Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. O1 replication journey – part 2: Surpassing
407 o1-preview through simple distillation, big progress or bitter lesson?, 2024.
- 408 [8] Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang,
409 and Mengnan Du. The impact of reasoning step length on large language models, 2024.
- 410 [9] Pravneet Kaur, Gautam Siddharth Kashyap, Ankit Kumar, Md Tabrez Nafis, Sandeep Kumar,
411 and Vikrant Shokeen. From text to transformation: A comprehensive review of large language
412 models’ versatility, 2024.
- 413 [10] Jack P.C. Kleijnen. *Design and Analysis of Simulation Experiments*, volume 230 of *International*
414 *Series in Operations Research & Management Science*. Springer International Publishing, Cham,
415 2015.
- 416 [11] Chaochao Lu, Chen Qian, Guodong Zheng, Hongxing Fan, Hongzhi Gao, Jie Zhang, Jing
417 Shao, Jingyi Deng, Jinlan Fu, Kexin Huang, Kunchang Li, Lijun Li, Limin Wang, Lu Sheng,
418 Meiqi Chen, Ming Zhang, Qibing Ren, Sirui Chen, Tao Gui, Wanli Ouyang, Yali Wang, Yan
419 Teng, Yaru Wang, Yi Wang, Yinan He, Yingchun Wang, Yixu Wang, Yongting Zhang, Yu Qiao,
420 Yujiong Shen, Yurong Mou, Yuxi Chen, Zaibin Zhang, Zhelun Shi, Zhenfei Yin, and Zhipin
421 Wang. From gpt-4 to gemini and beyond: Assessing the landscape of mllms on generalizability,
422 trustworthiness and causality through four modalities, 2024.
- 423 [12] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
424 policy optimization algorithms, 2017.
- 425 [13] Lewis Tunstall, Leandro Von Werra, and Thomas Wolf. *Natural Language Processing with*
426 *Transformers*. " O’Reilly Media, Inc.", 2022.
- 427 [14] Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don’t
428 always say what they think: Unfaithful explanations in chain-of-thought prompting. In *Proceed-*
429 *ings of the 37th International Conference on Neural Information Processing Systems*, NIPS ’23,
430 pages 74952–74965, Red Hook, NY, USA, 2023. Curran Associates Inc.
- 431 [15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
432 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st*
433 *International Conference on Neural Information Processing Systems*, NIPS’ 17, pages 6000–
434 6010, Red Hook, NY, USA, 2017. Curran Associates Inc.

- 435 [16] Jiaan Wang, Fandong Meng, Yunlong Liang, and Jie Zhou. Drt: Deep reasoning translation via
436 long chain-of-thought, 2025.
- 437 [17] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
438 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
439 models. In *Proceedings of the 36th International Conference on Neural Information Processing*
440 *Systems*, NIPS '22, pages 24824–24837, Red Hook, NY, USA, 2022. Curran Associates Inc.
- 441 [18] Yuxiang Zhang, Yuqi Yang, Jiangming Shu, Yuhang Wang, Jinlin Xiao, and Jitao Sang. Openrft:
442 Adapting reasoning foundation model for domain-specific tasks with reinforcement fine-tuning,
443 2024.
- 444 [19] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, and et al. Evaluation of openai o1:
445 Opportunities and challenges of agi, 2024.