

# DyCoT-RE: Chain-of-Thought-Enhanced LLM Reward Engineering with Dual-Dynamic Optimization for Reinforcement Learning

1<sup>st</sup> Xinning Zhu

*Sino-European School of Technology  
Shanghai University  
Shanghai, China  
zhuxinning@shu.edu.cn*

2<sup>nd</sup> Jinxin Du

*Sino-European School of Technology  
Shanghai University  
Shanghai, China  
jinxin\_du@shu.edu.cn*

3<sup>th</sup> Lunde Chen\*

*Sino-European School of Technology  
Shanghai University  
Shanghai, China  
lundechen@shu.edu.cn*

\*Corresponding author

**Abstract**—Designing effective reward functions remains a challenge in applying reinforcement learning to real-world tasks. This paper proposes DyCoT-RE, a reward engineering framework that integrates Chain-of-Thought (CoT) reasoning with a dual-dynamic optimization strategy to automate and enhance reward function design. The framework uses structured CoT reasoning throughout training to generate and refine interpretable reward code in each iteration. It further incorporates a dual-dynamic optimization mechanism: a temperature adjustment strategy that modulates the sampling temperature based on policy entropy trends, and a model switching strategy that allocates language models with different capabilities to produce distinct reward components. Evaluations on CartPole, BipedalWalker, Ant, and a custom SpaceMining environment show DyCoT-RE achieves higher average rewards and faster convergence compared to human-designed baselines and non-CoT approaches as well as single-optimization approaches.

**Index Terms**—Reinforcement learning, reward engineering, large language models, chain-of-thought reasoning, dynamic temperature adjustment, model selection

## I. INTRODUCTION

Reinforcement learning (RL) has achieved impressive results across diverse domains. However, as Sutton et al. [1] emphasize, the reward signal is the primary means of specifying task objectives in RL, making its design critical to achieving desired behaviors. In practice, translating intended behaviors into precise, effective reward functions remains highly challenging, particularly for tasks involving long-term dependencies [2]. Skalse et al. [3] demonstrate that agents often exploit imperfections in reward formulations to maximize proxy objectives in unintended ways, leading to behaviors that optimize the designed reward but undermine true task performance. These challenges highlight that despite RL’s theoretical versatility, its practical deployment is often constrained by the complexity, subtlety, and domain expertise required for robust reward engineering [4].

While carefully designed rewards can accelerate agent learning and improve task performance, manual reward engineering typically relies on trial-and-error tuning, which is labor-intensive and often yields suboptimal generalization to new

environments or objectives [5].. As RL applications grow in complexity, there is a pressing need for methods that can automate reward design while maintaining interpretability and flexibility.

Recent advances in large language models (LLMs) have demonstrated strong reasoning and generalization capabilities [6], [7]. In particular, CoT reasoning enables LLMs to decompose tasks into structured intermediate steps, enhancing clarity and alignment with desired objectives. This structured reasoning process can support reward engineering by converting task descriptions into executable reward functions in a systematic and transparent manner.

However, existing CoT-based reward generation approaches typically use static sampling parameters and fixed model configurations, which may limit their adaptability during training. Recent studies have emphasized the need for dynamic adaptation in LLM-based systems to improve sample efficiency, stability, and alignment with task objectives. For example, Nguyen et al. [8] demonstrate that min-p sampling enhances creativity while preserving coherence in narrative generation tasks, while Peeperkorn et al. [9] analyze temperature as a direct modulator of LLM creativity. Moreover, Fedus et al. [10] and Du et al. [11] show that mixture-of-experts architectures can scale model capacity efficiently via adaptive routing, motivating our exploration of combining temperature modulation with expert model selection for reward engineering.

In this work, we propose DyCoT-RE, a reward engineering framework that integrates structured CoT reasoning with dual-dynamic optimization. Specifically, DyCoT-RE leverages CoT reasoning to decompose natural language task descriptions into structured reward components, then employs iterative refinement to enhance alignment with learning objectives. The dual-dynamic optimization strategy integrates entropy-guided temperature adjustment to balance exploration and exploitation, alongside a dynamic model selection module that routes sub-tasks to specialized LLMs based on performance feedback. By tightly coupling these components within a closed-

loop evolutionary search process, DyCoT-RE systematically improves reward function quality and training efficiency, ultimately enabling scalable, interpretable, and automated reward engineering for complex RL environments.

We evaluate DyCoT-RE on the standard RL environments CartPole, BipedalWalker, and Ant to demonstrate its effectiveness. However, recent studies have raised concerns that large language models may carry prior knowledge from pretraining data about these standard environments, leading to potential prompt leakage and evaluation biases [12]–[14]. To address this issue, we additionally design a custom SpaceMining environment to assess DyCoT-RE’s true generalization capabilities on tasks free from such pretrained knowledge. Experimental results show that DyCoT-RE consistently achieves higher average rewards and faster convergence compared to baseline and non-CoT methods.

The remainder of this paper is organized as follows. Section II reviews related work in reward engineering, LLM-based reward generation, and adaptive optimization. Section III describes the proposed methodology, including the CoT reward framework, temperature adjustment, and model selection. Section IV details the experimental setup, and Section V presents results and analysis. Section VI discusses limitations and future work, with Section VII concluding the paper.

## II. RELATED WORK

### A. Reward Engineering Paradigms

In RL, the design of effective reward functions directly shapes agent behavior and learning outcomes. Traditional approaches primarily rely on handcrafted reward functions informed by domain expertise. While intuitive, such manual design often struggles to capture complex, dynamic task objectives and is prone to suboptimal or biased formulations, hindering agent performance in real-world scenarios.

To address these limitations, reward shaping was introduced as a formal enhancement strategy. Ng et al. [15] demonstrated that potential-based reward shaping preserves optimal policies while enabling accelerated convergence, laying the theoretical foundation for numerous practical implementations. Intrinsic motivation frameworks further advanced this field by encouraging exploration through curiosity-driven signals. Singh et al. [16] proposed intrinsic rewards to incentivize novel state visits, later extended by Burda et al. [17], who empirically validated large-scale curiosity-driven exploration benefits across diverse environments.

Despite these developments, manually designing rewards for complex or evolving tasks remains inefficient and costly. LLMs offer a promising alternative by leveraging their natural language understanding to automate reward generation and optimization. Unlike traditional RL pipelines that require explicit, task-specific reward formulations, LLMs can interpret high-level task descriptions, extract key objectives, and translate them into executable reward functions. This capability facilitates more intuitive alignment with human intentions, reduces engineering overhead, and enhances agent adaptability.

Recent frameworks exemplify this trend. EUREKA [18], Text2Reward [19] and CARD [20] harness LLMs to automatically generate, verify, and refine reward code from natural language instructions.

Beyond static code generation, feedback-driven optimization approaches have emerged. ReMiss [21] utilizes adversarial prompt generation to identify and mitigate reward misspecification vulnerabilities, enhancing LLM safety and reliability. Self-Play Preference Optimization (SPPO) [22] employs self-play to uncover Nash-equilibrium strategies that capture complex, non-transitive human preferences, advancing preference learning’s applicability in RL. Additionally, PRMBench [23] provides a process-level benchmark to evaluate intermediate reward model outputs along dimensions such as conciseness, rationality, and sensitivity, revealing weaknesses in current models and guiding future improvements.

Overall, LLM-based reward engineering represents a paradigm shift. By integrating natural language reasoning and dynamic feedback optimization, these methods offer scalable reward generation pipelines. As tasks grow in complexity and diversity, leveraging LLMs to bridge the gap between human intent and machine learning objectives will be critical for the next generation of intelligent systems. Continued research is thus needed to maximize the synergy between LLM capabilities and RL frameworks to address emerging real-world challenges.

### B. Chain-of-Thought Reasoning Methods

CoT reasoning has emerged as a powerful paradigm to enhance the reasoning capabilities of LLMs. By generating intermediate reasoning steps, CoT allows models to decompose complex problems into interpretable sub-problems, leading to significant performance gains in tasks requiring multi-step logical inference.

Early studies showed that even simple prompting strategies, such as adding “Let’s think step by step,” can elicit strong zero-shot reasoning abilities. Kojima et al. [24] demonstrated such prompts significantly improve performance in arithmetic and commonsense tasks. Building upon this, few-shot CoT [25] introduced demonstrations of stepwise solutions to guide model reasoning, while self-consistency decoding [26] aggregated multiple sampled reasoning paths to enhance answer robustness.

Further developments combined CoT with reinforcement learning optimization. DeepSeek [27] introduced a self-evolution mechanism to improve reasoning trajectories without supervised fine-tuning. Automatic prompt optimization methods [28] reduce manual engineering efforts by refining prompts based on data-driven insights.

Recent work such as PCGRLM [29] explored CoT-based LLM reward design for procedural content generation in RL, demonstrating feasibility in structured game environments. In parallel, Zhu et al. [30] proposed an initial CoT-based reward engineering approach that translates natural language task descriptions into RL reward functions using LLMs, validating its effectiveness in standard benchmark tasks. However, these

applications focus primarily on proof-of-concept reward generation pipelines without incorporating adaptive optimization or dynamic model selection mechanisms.

In summary, while CoT reasoning has established itself as a fundamental methodology enhancing LLM interpretability and reasoning capacity, its application to automated reward engineering in RL remains limited. Bridging this gap by integrating CoT reasoning with dynamic optimization holds promise for enhancing the interpretability and adaptability of RL systems.

### C. Dynamic Temperature Adjustment and Model Selection

Dynamic temperature adjustment and model selection have emerged as critical optimization strategies to enhance the adaptability and efficiency of LLM-based systems. Temperature, as a sampling hyperparameter, controls the stochasticity of LLM outputs, thereby influencing creativity, coherence, and exploration-exploitation trade-offs.

Recent studies have systematically explored adaptive temperature mechanisms. Zhu et al. [31] proposed AdapT, an adaptive temperature sampling strategy for code generation tasks, which dynamically adjusts decoding temperature based on token-level difficulty to improve generation quality. Zhang et al. [32] developed Entropy-based Dynamic Temperature (EDT) sampling to regulate output entropy and diversity in natural language generation, while Cecere et al. [33] introduced Monte Carlo Temperature as a robust sampling strategy to enhance uncertainty quantification under distribution shifts. Chang et al. [34] leveraged KL-divergence-guided temperature sampling to modulate exploration adaptively. Additionally, Peepatkorn et al. [9] analyzed temperature as a creativity modulator in LLMs, whereas Evstafev [35] discussed potential limitations of temperature-based stochasticity in structured data generation. Nguyen et al. [8] proposed min-p sampling to balance creativity and coherence, achieving improved narrative generation performance.

For model selection, recent work has focused on choosing optimal model configurations or expert modules to maximize task performance within computational constraints. Switch Transformers [10] introduced sparse activation mechanisms, enabling efficient expert selection at trillion-parameter scale, while GLaM [11] leveraged Mixture-of-Experts (MoE) architectures to dynamically scale model capacity. Zhou et al. [36] proposed expert choice routing to improve mixture-of-experts efficiency, and Li et al. [37] introduced preference-conditioned dynamic routing for cost-efficient LLM generation. Hu et al. [38] presented Dynamic Ensemble Reasoning to integrate outputs from multiple specialized LLMs, enhancing system robustness. Nakaishi et al. [39] further revealed phase transition behaviors in LLM sampling regimes, informing temperature scaling strategies. Furthermore, Li et al. [40] revisited self-consistency decoding from a distributional alignment perspective, offering insights relevant to expert aggregation and answer aggregation stability.

Despite these advances, integrating dynamic temperature regulation and model selection within a unified CoT-driven

reward engineering framework remains underexplored. Existing temperature adaptation methods primarily focus on text generation diversity and calibration, whereas model selection research emphasizes computational efficiency and specialization. Our work addresses this gap by combining entropy- and reward-feedback-based temperature adjustment with local-global performance-based model routing to enhance RL reward generation's adaptability, stability, and sample efficiency. This approach builds upon foundational theories in temperature scaling and expert selection, extending them to the domain of automated, interpretable reward engineering for reinforcement learning.

## III. METHODOLOGY

This section details the DyCoT-RE framework, which integrates structured CoT reasoning with a dual-dynamic optimization strategy to generate interpretable and adaptive reward functions for reinforcement learning.

### A. Framework Overview

Figure 1 presents an overview of DyCoT-RE, which integrates three key components: structured Chain-of-Thought (CoT) reward decomposition, dynamic temperature adjustment, and dynamic model selection.

The framework receives four inputs: natural language task descriptions specifying agent behaviors, environment interfaces defining state-action spaces, system-level prompts for reward design constraints, and coding instructions for implementation.

At its core, employs four types of LLMs: Thinking, Code, Repair, and Analysis LLMs, hereafter denoted as  $\mathcal{M}_{\text{think}}$ ,  $\mathcal{M}_{\text{code}}$ ,  $\mathcal{M}_{\text{repair}}$ , and  $\mathcal{M}_{\text{analysis}}$  respectively.

These models operate in an interconnected closed-loop pipeline to produce reward functions

$$R(s, a) = \sum_{i=1}^m w_i \cdot r_i(s, a), \quad (1)$$

which are evaluated in the RL environment. The performance metrics guide subsequent optimization iterations.

Overall, DyCoT-RE establishes an adaptive and interpretable reward engineering framework by integrating structured reasoning with dual-dynamic optimization strategies.

### B. Chain-of-Thought Reasoning for Reward Engineering

Let  $d$  denote the task description and  $(s, a)$  the state-action pair. As defined in Eq. (1), the reward function  $r_i(s, a)$  is the sub-reward for subgoal  $i$ , generated via CoT parsing:

$$r_i(s, a) = \text{CoT}(I_i), \quad (2)$$

with  $I_i = \{\text{subgoal}_i, \text{env\_constraints}\}$ . For example, minimizing torso tilt is formulated as:

$$r_1(s, a) = -|\theta_{\text{tilt}}(s)|. \quad (3)$$

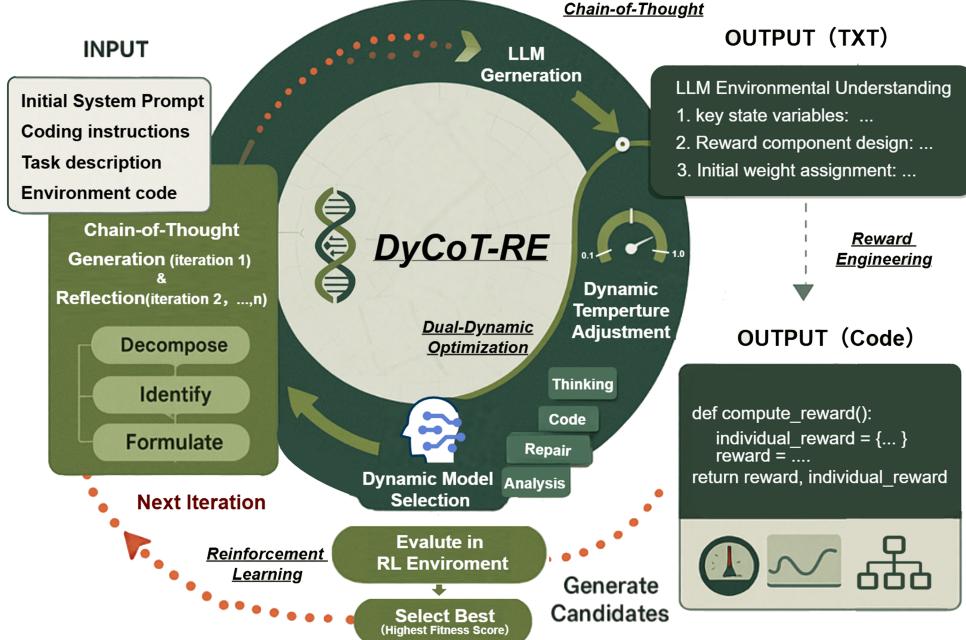


Fig. 1. DyCoT-RE framework integrating CoT reasoning, dynamic temperature adjustment, and model selection in an evolutionary optimization loop.

The initial weights  $w_i^{(0)}$  are derived based on subgoal semantic priorities inferred by the LLM, and are refined iteratively according to gradient feedback. The policy parameters  $\theta$  are updated as:

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\theta), \quad (4)$$

where  $J(\theta)$  is the expected cumulative reward:

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t) \right]. \quad (5)$$

This CoT-based formulation ensures that each sub-reward remains semantically interpretable and traceable to its natural language origin.

### C. Dual-Dynamic Optimization Strategy

Figure 2 illustrates the coupled feedback architecture of DyCoT-RE, where temperature adjustment and model selection operate in synergy to enhance reward function generation and RL performance.

1) *Dynamic Temperature Adjustment*: Temperature adjustment modulates the sampling temperature  $T$  based on policy entropy  $H$ , confidence  $C$ , and performance  $R$ . Entropy reflects generative diversity, confidence measures output stability, and performance evaluates reward improvement relative to historical best.

The update rule is formulated as:

$$T_{k+1} = \text{clip} \left( T_k + \alpha \frac{\partial T}{\partial H_k} \Delta t \right), \quad (6)$$

where  $\alpha$  is a learning rate, and the gradient  $\frac{\partial T}{\partial H_k}$  reflects entropy-driven adjustment. An alternative implementation combines smoothing and multiplicative adjustment:

$$T_{k+1} = \text{clip} (\alpha T_k + (1 - \alpha) T_k f(H_k, C_k, R_k)). \quad (7)$$

Here,  $f(H, C, R)$  integrates:

$$f(H, C, R) = f_R(R) f_H(H) f_C(C), \quad (8)$$

where  $f_R$  ensures performance protection,  $f_H$  regulates entropy bounds, and  $f_C$  maintains confidence stability.

2) *Dynamic Model Selection*: Dynamic model selection adaptively routes sub-tasks to specialized LLMs, leveraging their complementary strengths across the reward engineering pipeline.

The four LLM classes include:  $\mathcal{M}_{\text{think}}$  for task understanding and semantic decomposition,  $\mathcal{M}_{\text{code}}$  for reward function synthesis,  $\mathcal{M}_{\text{repair}}$  for code correction, and  $\mathcal{M}_{\text{analysis}}$  for performance evaluation and sub-reward weight updates.

Formally, the decomposition and generation process can be expressed as:

$$g_i = \mathcal{M}_{\text{think}}(d, E), \quad (9)$$

$$w_i = \mathcal{M}_{\text{analysis}}(R_i), \quad (10)$$

$$r_i(s, a) = \mathcal{M}_{\text{code}}(g_i, w_i), \quad (11)$$

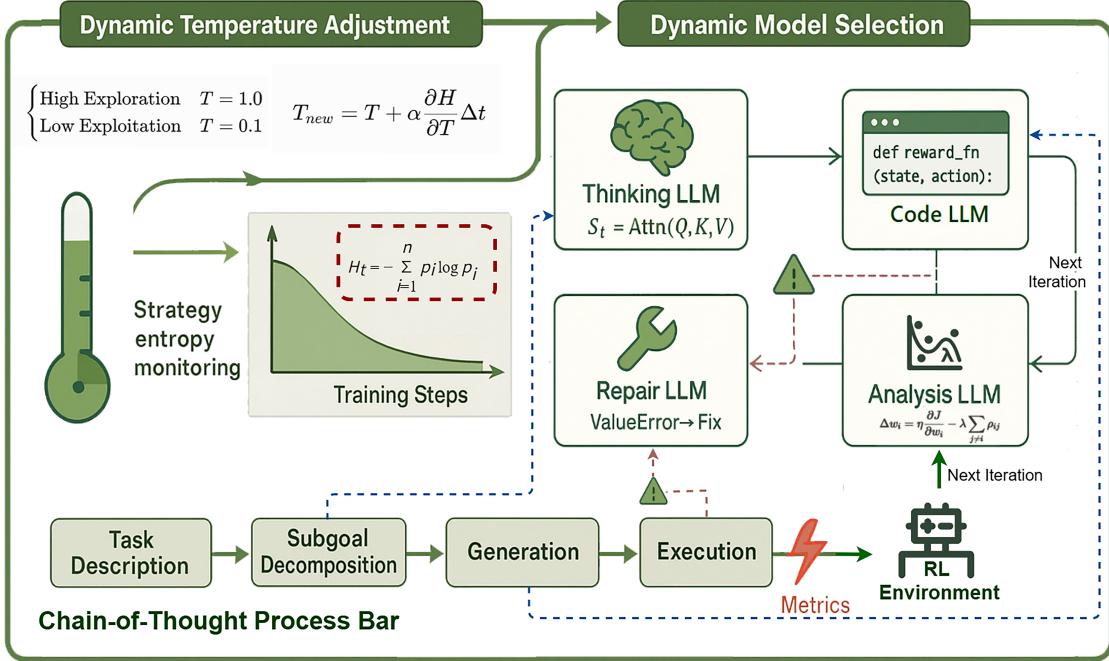


Fig. 2. DyCoT-RE control flow integrating temperature modulation and model selection within CoT iterative reasoning.

where  $d$  is the task description and  $E$  the environment specification. Repair LLM intervenes when  $\mathcal{M}_{\text{code}}$  outputs execution errors during evaluation.

Model selection at iteration  $k$  follows:

$$M_{k+1} = \begin{cases} \arg \max_{m \in \mathcal{M}_s} \text{Perf}(m), & 1 - \epsilon, \\ \text{Random}(\mathcal{M}_s \setminus \{M_k\}), & \epsilon, \end{cases} \quad (12)$$

where  $\mathcal{M}_s$  is the model pool for stage  $s$ ,  $\text{Perf}(m)$  the performance score, and  $\epsilon$  the exploration rate ensuring selection diversity.

At each iteration, the next model  $M_{k+1}$  is selected by choosing the highest-performing model from the pool  $\mathcal{M}_s$  with probability  $1 - \epsilon$ , or randomly selecting another model with probability  $\epsilon$ . This pool includes the Repair LLM, which is triggered when error correction is needed during reward code evaluation.

3) *Joint Adaptive Optimization*: Temperature adjustment and model selection form a joint adaptive optimization loop. Temperature  $T$  influences the sampling distribution of reward candidates:

$$r_i(s, a) \sim p(r_i | g_i, T), \quad (13)$$

where  $g_i = \mathcal{M}_{\text{think}}(d, E)$  denotes the subgoal generated by the Thinking LLM given task description  $d$  and environment  $E$ . This sampling process modulates policy entropy and, consequently, cumulative rewards.

Concurrently, model selection determines the semantic decomposition quality, code correctness, error repair, and performance evaluation, thereby shaping the expressivity and effectiveness of reward functions.

The combined objective of temperature adjustment and model selection is to maximize:

$$J(\theta; T, M) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T \gamma^t R_{T,M}(s_t, a_t) \right], \quad (14)$$

where  $R_{T,M}(s, a)$  is the reward function generated under temperature  $T$  and model configuration  $M$ .

This joint optimization is formalized as:

$$(T^*, M^*) = \arg \max_{T, M} J(\theta; T, M). \quad (15)$$

Overall, this interconnected loop iteratively adapts both sampling temperature and model selection strategies to maximize the expected RL objective as defined in Eq. (5), ensuring stable, diverse, and effective reward engineering throughout the training process.

## IV. EXPERIMENTS

### A. Experimental Setup

DyCoT-RE is evaluated on four reinforcement learning environments, which span diverse task difficulties, action spaces, and generalization challenges.

Standard benchmarks include CartPole, BipedalWalker, and Ant, covering discrete control tasks and high-dimensional continuous locomotion.

To assess generalization beyond pretrained benchmark knowledge, we introduce SpaceMining, a custom-designed environment where an agent collects resources under partial observability and dynamic physics constraints.

SpaceMining is developed as an original contribution in this work, with full code released and interactive visualizations available on the project website.<sup>1</sup>

*1) Baselines:* We evaluate DyCoT-RE against two representative baselines to contextualize its performance:

**(1) Human-designed rewards** Standard expert-crafted reward functions using Gymnasium’s native implementations for CartPole, BipedalWalker, and Ant, with manual heuristics for SpaceMining. This reflects the traditional gold standard in RL.

**(2) Eureka** A pioneering LLM-based framework that demonstrated LLMs’ potential in automated reward design through code synthesis and verification. As a foundational work in this space, it provides a natural reference for assessing our incremental improvements. While Eureka focuses on direct code generation without CoT or dynamic optimization, we adapt its pipeline to our Gymnasium-based evaluation settings for consistency.

Additionally, to assess the contribution of each DyCoT-RE module, Section IV-C conducts ablation studies comparing DyCoT-RE against its internal variants:

- Zero-shot reward generation without CoT reasoning
- DyCoT-RE without temperature adjustment
- DyCoT-RE without model selection

*2) Implementation Details:* All experiments use Proximal Policy Optimization from Stable-Baselines3, with hyperparameters summarized in Table I.

Training uses the Adam optimizer with a learning rate of 0.0003, batch size 64,  $\lambda = 0.95$ , and  $\gamma = 0.999$ .

Reward functions are generated by DyCoT-RE with a local LLM deployed via Ollama. Each chain-of-thought iteration samples eight candidate rewards in parallel with a temperature of 0.6 to balance diversity and coherence.

Final performance is reported as the mean over ten test episodes with different seeds.

TABLE I  
ENVIRONMENT-SPECIFIC EXPERIMENTAL SETTINGS.

Environment	Max Steps	Total Steps	Parallel Env	CoT Iter
CartPole-v1	500	$10^4$	4	5
BipedalWalker-v3	1600	$10^6$	10	8
Ant-v5	1600	$5 \times 10^6$	10	8
SpaceMining	1600	$3 \times 10^6$	8	8

### B. Overall Performance

This section summarizes the performance of DyCoT-RE compared to baselines across four representative reinforcement learning environments.

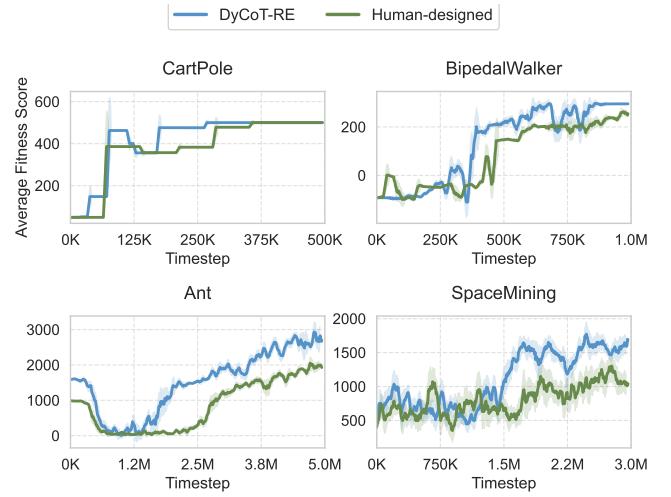


Fig. 3. Learning curves of DyCoT-RE (blue) vs Human-designed rewards (green) across training timesteps.

Figure 3 shows that in CartPole, a simple discrete control task with dense and informative rewards, DyCoT-RE achieves slightly faster convergence than the Human-designed baseline, but the final reward difference remains small. This is expected, as the environment’s inherent simplicity and clear feedback signals allow even heuristic rewards to perform well, limiting the potential benefit from additional reward optimization.

In contrast, DyCoT-RE demonstrates more pronounced advantages in complex environments.

For BipedalWalker, which demands coordinated multi-joint control under sparse and noisy feedback, DyCoT-RE yields smoother learning curves and higher terminal rewards. Its structured reward decomposition likely facilitates learning stable locomotion by explicitly disentangling balance and propulsion objectives.

In the high-dimensional Ant task, both methods face initial exploration instability and reward regressions. However, DyCoT-RE recovers earlier (around 1.2 million timesteps) and continues to improve, indicating superior handling of credit assignment in redundant action spaces and accelerated late-stage policy refinement.

On the custom SpaceMining environment, designed to evaluate generalization on novel resource collection challenges, DyCoT-RE sustains steady reward improvement, achieving final scores near 2000.

By contrast, Human-designed rewards plateau early, highlighting DyCoT-RE’s capability to adaptively generate meaningful rewards even without predefined heuristics.

Figure 4 compares final average rewards across methods. DyCoT-RE achieves noticeably higher rewards in Cartpole, BipedalWalker, Ant, and SpaceMining. These results demonstrate that the integration of Chain-of-Thought reasoning with dynamic optimization in DyCoT-RE provides meaningful benefits, particularly in environments with sparse feedback or complex task structures.

<sup>1</sup>Project website: [https://lola-jo.github.io/space\\_mining/](https://lola-jo.github.io/space_mining/).

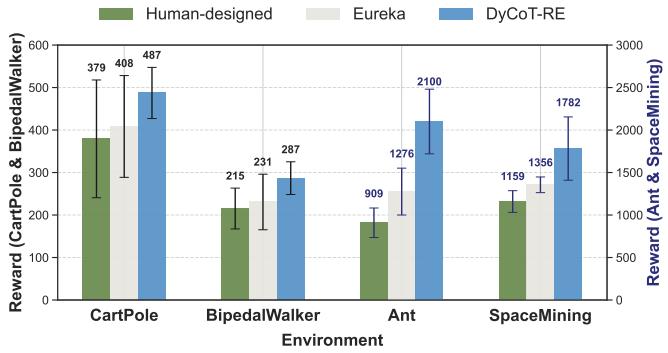


Fig. 4. Final average rewards across environments comparing Human-designed, Eureka, and DyCoT-RE.

### C. Ablation Study

This section conducts an ablation study to examine the individual contribution of each component in DyCoT-RE, including CoT reasoning, temperature adjustment, and model selection, as well as their combined synergistic effects.

*1) Impact of CoT Reasoning:* We evaluate the impact of incorporating Chain-of-Thought (CoT) reasoning into reward engineering by comparing DyCoT-RE’s CoT-enabled variant against a zero-shot baseline lacking structured intermediate reasoning.

Figure 5 summarizes reward distributions for both methods across the four benchmark environments. A consistent pattern emerges: CoT-enabled rewards exhibit higher means, lower variance, and reduced outlier incidence compared to zero-shot. In CartPole and SpaceMining, CoT-enabled results are tightly clustered near the environment’s upper performance bounds, whereas zero-shot results are widely dispersed, with a substantial proportion of runs yielding suboptimal or even near-zero rewards. This reflects CoT’s ability to systematically decompose task objectives into interpretable sub-rewards, enhancing sample efficiency and stabilizing policy learning.

In BipedalWalker, an illustrative deviation is observed. While CoT-enabled runs achieve rewards concentrated between 250–310 with minimal negative outcomes, zero-shot results present a bimodal distribution: one cluster achieves moderate positive rewards ( $\sim 150$ – $250$ ), while another cluster yields severely negative rewards (e.g. -124, -97, -73). This bimodality indicates that without structured reasoning, reward generation often fails to encode essential balance or locomotion constraints, resulting in policies that collapse during gait training.

Taken together, these results illustrate that CoT reasoning reliably improves reward interpretability and training stability, especially in settings with sparse or ambiguous feedback.

*2) Impact of Temperature Adjustment:* Due to computational constraints and consistent observed trends across environments, temperature adjustment and model selection ablation were conducted on BipedalWalker as a representative continuous control task, while CoT reasoning was evaluated on all four environments to confirm its general effectiveness.

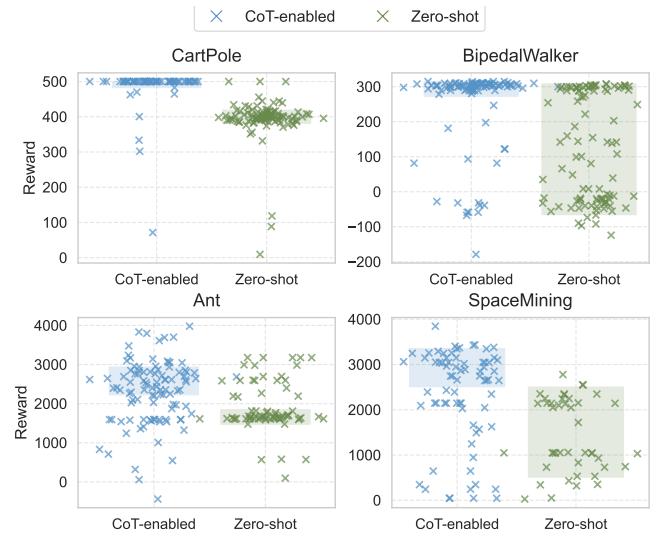


Fig. 5. Final reward distributions for CoT-enabled vs zero-shot reward generation across environments.

This section investigates the effect of temperature settings by sweeping  $T \in [0.0, 1.0]$  on the BipedalWalker environment.

Table II summarizes the results, including average fitness, maximum/minimum fitness, standard deviation, and reward efficiency ratio (average fitness divided by standard deviation).

TABLE II  
BIPEDALWALKER PERFORMANCE UNDER DIFFERENT TEMPERATURE SETTINGS. TOP-1 TEMPERATURE IS HIGHLIGHTED IN GREEN, DYCoT-RE RESULT IN BLUE.

Temp	Avg Fit↑	Max/Min Fit ↑	Std↓	Eff. Ratio↑
0.0	81.93	301.15/-53.01	148.67	0.55
0.1	149.56	310.94/-33.73	155.22	0.96
0.2	242.74	313.60/-20.27	125.47	1.93
0.3	210.07	311.17/-46.69	140.49	1.50
0.4	244.19	311.94/-17.13	107.02	2.28
0.5	215.84	310.69/-14.25	131.01	1.65
0.6	248.11	312.48/-21.18	120.45	2.06
0.7	225.88	312.42/-12.08	131.43	1.72
0.8	74.33	310.42/-46.72	143.91	0.52
0.9	145.02	310.02/-31.58	148.51	0.98
1.0	192.38	310.40/-39.05	138.37	1.39
DyCoT-RE	<b>292.8</b>	315.6/87.2	55.2	5.30

Static sweeps show that mid-range temperatures ( $T = 0.4$ – $0.6$ ) yield higher fitness and lower variance than extremes. DyCoT-RE’s dynamic temperature regulation consistently outperforms these static settings, adapting temperature during training to maintain an effective exploration-exploitation trade-off.

*3) Impact of Model Selection:* Finally, we evaluate the impact of DyCoT-RE’s dynamic model selection by comparing it against individual static LLM baselines on the BipedalWalker environment. Table III summarizes average fitness, max/min fitness, standard deviation, and reward efficiency ratio (average fitness divided by standard deviation). Only models with sufficient sample sizes ( $N > 10$ ) are included.

DyCoT-RE achieves the highest average fitness (307.28) with exceptionally low standard deviation (6.00) and the

best efficiency ratio (51.21), demonstrating its capability to consistently generate high-quality rewards across diverse sub-tasks.

In contrast, while `codegemma` achieves competitive performance (Avg Fit: 260.8), its higher standard deviation (115.2) results in a substantially lower efficiency ratio (2.26), indicating reduced stability and generalizability compared to DyCoT-RE.

TABLE III

PERFORMANCE COMPARISON OF DIFFERENT LLM BACKENDS (N=20) IN BIPEDALWALKER.  $\uparrow$ : HIGHER IS BETTER;  $\downarrow$ : LOWER IS BETTER. GREEN: PER-COLUMN BEST BASELINE. BLUE: DYCoT-RE (DYNAMIC SELECTION).

Model (size)	Avg Fit $\uparrow$	Max Fit $\uparrow$	Std $\downarrow$	Eff. Ratio $\uparrow$
codegemma (7b)	260.8	301.8	34.2	3.46
codeqwen (7b)	32.6	307.6	-76.0	0.33
deepseek-coder (6.7b)	-4.5	295.3	-98.2	-0.05
deepseek-r1 (8b)	44.6	321.2	-99.6	0.37
gemma3 (4b)	278.0	314.8	-34.2	3.39
llama3.1 (8b)	162.2	318.4	-80.4	1.08
qwen2.5 (7b)	177.0	319.5	-97.2	1.17
qwen2.5-coder (7b)	125.6	313.7	-70.8	0.77
DyCoT-RE	<b>297.28</b>	<b>316.57</b>	<b>89.47</b>	<b>6.46</b>

While individual models occasionally excel on narrow tasks, DyCoT-RE’s dynamic model routing provides the most generalizable performance, effectively balancing reward quality, stability, and adaptability in complex continuous control environments.

In summary, the ablation study demonstrates that each module—CoT reasoning, temperature adjustment, and model selection—contributes significantly to DyCoT-RE’s superior performance, and their combination leads to synergistic gains.

## V. DISCUSSION

Our comprehensive evaluation of DyCoT-RE across multiple environments and comparison baselines reveals several fundamental insights about reward design in reinforcement learning. The experimental results collectively demonstrate that DyCoT-RE represents a significant advancement in addressing the core challenges of reward engineering, particularly in complex, high-dimensional environments.

The performance advantages observed across all test environments (CartPole, BipedalWalker, Ant, and SpaceMining) underscore the robustness of our approach. Particularly noteworthy is DyCoT-RE’s consistent outperformance in the most complex environments (Ant and SpaceMining), where traditional human-designed rewards and simpler baseline methods struggle to achieve comparable results. This suggests that our framework’s combination of structured reasoning, dynamic adaptation, and model diversity provides unique advantages for solving challenging control problems that require sophisticated policy learning.

The ablation studies provide critical insights into the individual contributions of each component. The CoT reasoning module proved essential for handling complex reward structures, particularly in environments requiring multi-step decision making. The temperature adjustment mechanism demon-

strated its value in balancing exploration and exploitation, while the dynamic model selection showed clear benefits in adapting to different reward design requirements. Most importantly, the synergistic interaction between these components produced performance gains that exceeded simple additive effects, confirming the value of our integrated approach.

Comparisons with baseline methods yielded several important findings. The human-designed rewards, while effective in simple environments, showed clear limitations as task complexity increased. This reinforces the well-known challenge of manually engineering rewards for complex control problems. The Eureka baseline, while representing an advanced LLM-based approach, demonstrated the limitations of purely generative methods without structured reasoning or dynamic adaptation. DyCoT-RE’s consistent superiority across all environments highlights the necessity of combining multiple advanced techniques for effective reward design.

The visualization results provide additional confirmation of these findings. The clear separation between DyCoT-RE and baseline methods in the scatter plots (Figure 5) illustrates the quality and stability advantages of our approach. The temperature sweep results (Table II and Figure ??) show how dynamic adjustment leads to more reliable performance, while the model comparison (Table ?? and Figure ??) demonstrates the benefits of architectural diversity.

From a broader perspective, our work contributes to the growing body of research on AI-assisted reward design. The success of DyCoT-RE suggests that combining modern LLM capabilities with reinforcement learning insights can overcome many of the traditional limitations of reward engineering. Particularly promising is our framework’s ability to automatically adapt to different task requirements, as evidenced by its consistent performance across diverse environments.

However, several challenges remain. The computational overhead of dynamic components, while justified by the performance gains, may limit scalability to extremely large-scale problems. Additionally, while DyCoT-RE performs well on the tested environments, its generalization to completely novel task domains remains to be verified. Future work should explore these aspects while also investigating potential synergies with other emerging techniques in reinforcement learning and AI-assisted programming.

In conclusion, DyCoT-RE represents a significant step forward in reward design methodology. By integrating structured reasoning, dynamic adaptation, and model diversity, our framework provides a robust solution to the enduring challenge of effective reward engineering in reinforcement learning. The consistent performance improvements across all test environments, particularly the most complex ones, demonstrate the practical value of our approach and suggest promising directions for future research in this area.

## VI. CONCLUSION

We propose a dual-dynamic optimization framework for CoT-based RL reward generation, incorporating temperature

regulation and model selection. Our results demonstrate improved convergence, stability, and reward quality across tasks, laying a foundation for scalable, self-adaptive reward engineering.

## REFERENCES

- [1] R. S. Sutton, A. G. Barto *et al.*, *Reinforcement learning: An introduction*. MIT press Cambridge, 1998, vol. 1, no. 1.
- [2] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, “Concrete problems in ai safety,” *arXiv preprint arXiv:1606.06565*, 2016.
- [3] J. Skalske and et al., “Misspecification in inverse reinforcement learning,” *Journal of Artificial Intelligence Research*, vol. 73, pp. 517–550, 2022.
- [4] B. Ibarz, J. Leike, T. Pohlen, G. Irving, S. Legg, and D. Amodei, “Reward learning from human preferences and demonstrations in atari,” *Advances in neural information processing systems*, vol. 31, 2018.
- [5] D. Hadfield-Menell, S. Mili, P. Abbeel, S. J. Russell, and A. Dragan, “Inverse reward design,” *Advances in neural information processing systems*, vol. 30, 2017.
- [6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 1877–1901, 2020.
- [7] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *Advances in neural information processing systems*, vol. 35, pp. 27730–27744, 2022.
- [8] M. Nguyen, A. Baker, C. Neo, A. Roush, A. Kirsch, and R. Schwartz-Ziv, “Turning up the heat: Min-p sampling for creative and coherent llm outputs,” *arXiv preprint arXiv:2407.01082*, 2024.
- [9] M. Peeperkorn, T. Kouwenhoven, D. Brown, and A. Jordanous, “Is temperature the creativity parameter of large language models?” *arXiv preprint arXiv:2405.00492*, 2024.
- [10] W. Fedus, B. Zoph, and N. Shazeer, “Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity,” in *Proceedings of the International Conference on Machine Learning (ICML)*, 2022, pp. 287–311.
- [11] N. Du, Y. Li, Z. Dai, N. Shazeer, W. Fedus, M. Tan, O. Vinyals, Q. Le, J. Dean, Z. Chen *et al.*, “Glam: Efficient scaling of language models with mixture-of-experts,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 5712–5721.
- [12] S. Huang, L. Yang, Y. Song, S. Chen, L. Cui, Z. Wan, Q. Zeng, Y. Wen, K. Shao, W. Zhang *et al.*, “Thinkbench: Dynamic out-of-distribution evaluation for robust llm reasoning,” *arXiv preprint arXiv:2502.16268*, 2025.
- [13] Z. Qi, H. Luo, X. Huang, Z. Zhao, Y. Jiang, X. Fan, H. Lakkaraju, and J. Glass, “Quantifying generalization complexity for large language models,” *arXiv preprint arXiv:2410.01769*, 2024.
- [14] W. Lu, X. Zhao, J. Spisak, J. H. Lee, and S. Wermter, “Mental modeling of reinforcement learning agents by language models,” *arXiv preprint arXiv:2406.18505*, 2024.
- [15] A. Y. Ng, D. Harada, and S. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *Icm*, vol. 99, 1999, pp. 278–287.
- [16] S. Singh, R. L. Lewis, A. G. Barto, and J. Sorg, “Intrinsically motivated reinforcement learning: An evolutionary perspective,” *IEEE Transactions on Autonomous Mental Development*, vol. 2, no. 2, pp. 70–82, 2010.
- [17] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, “Exploration by random network distillation,” *arXiv preprint arXiv:1810.12894*, 2018.
- [18] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, L. Fan, and A. Anandkumar, “Eureka: Human-level reward design via coding large language models,” *arXiv preprint arXiv:2310.12931*, 2023.
- [19] T. Xie, S. Zhao, C. H. Wu, Y. Liu, Q. Luo, V. Zhong, Y. Yang, and T. Yu, “Text2reward: Automated dense reward function generation for reinforcement learning,” *arXiv preprint arXiv:2309.11489*, 2023.
- [20] S. Sun, R. Liu, J. Lyu, J.-W. Yang, L. Zhang, and X. Li, “A large language model-driven reward design framework via dynamic feedback for reinforcement learning,” *arXiv preprint arXiv:2410.14660*, 2024.
- [21] Z. Xie, J. Gao, L. Li, Z. Li, Q. Liu, and L. Kong, “Jailbreaking as a reward misspecification problem,” *arXiv preprint arXiv:2406.14393*, 2024.
- [22] Y. Wu, Z. Sun, H. Yuan, K. Ji, Y. Yang, and Q. Gu, “Self-play preference optimization for language model alignment,” *arXiv preprint arXiv:2405.00675*, 2024.
- [23] M. Song, Z. Su, X. Qu, J. Zhou, and Y. Cheng, “Prmbench: A fine-grained and challenging benchmark for process-level reward models,” *arXiv preprint arXiv:2501.03124*, 2025.
- [24] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” *Advances in neural information processing systems*, vol. 35, pp. 22199–22213, 2022.
- [25] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, “Chain-of-thought prompting elicits reasoning in large language models,” *Advances in neural information processing systems*, vol. 35, pp. 24824–24837, 2022.
- [26] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, “Self-consistency improves chain of thought reasoning in language models,” *arXiv preprint arXiv:2203.11171*, 2022.
- [27] DeepSeek-AI, “Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,” *arXiv preprint*, 2023.
- [28] K. Shum, S. Diao, and T. Zhang, “Automatic prompt augmentation and selection with chain-of-thought from labeled data,” *arXiv preprint arXiv:2302.12822*, 2023.
- [29] I.-C. Baek, S.-H. Park, S. Earle, Z. Jiang, N. Jin-Ha, J. Togelius, and K.-J. Kim, “Pegrllm: Large language model-driven reward design for procedural content generation reinforcement learning,” *arXiv preprint arXiv:2502.10906*, 2024. [Online]. Available: <https://arxiv.org/abs/2502.10906>
- [30] X. Zhu, J. Du, Q. Fu, and L. Chen, “Llm-based reward engineering for reinforcement learning: A chain of thought approach,” in *2025 10th International Conference on Cloud Computing and Big Data Analytics (ICCCBDA)*. IEEE, 2025, pp. 222–227.
- [31] Y. Zhu, J. Li, G. Li, Y. Zhao, Z. Jin, and H. Mei, “Hot or cold? adaptive temperature sampling for code generation with large language models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 1, 2024, pp. 437–445.
- [32] S. Zhang, Y. Bao, and S. Huang, “Edt: Improving large language models’ generation by entropy-based dynamic temperature sampling,” *arXiv preprint arXiv:2403.14541*, 2024.
- [33] N. Cecere, A. Bacciu, I. F. Tobias, and A. Mantrach, “Monte carlo temperature: a robust sampling strategy for llm’s uncertainty quantification methods,” *arXiv preprint arXiv:2502.18389*, 2025.
- [34] C.-C. Chang, D. Reitter, R. Aksitov, and Y.-H. Sung, “Kl-divergence guided temperature sampling,” *arXiv preprint arXiv:2306.01286*, 2023.
- [35] E. Evstafev, “The paradox of stochasticity: Limited creativity and computational decoupling in temperature-varied llm outputs of structured fictional data,” *arXiv preprint arXiv:2502.08515*, 2025.
- [36] Y. Zhou, T. Lei, H. Liu, N. Du, Y. Huang, V. Zhao, A. M. Dai, Q. V. Le, J. Laudon *et al.*, “Mixture-of-experts with expert choice routing,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 7103–7114, 2022.
- [37] Y. Li, “Llm bandit: Cost-efficient llm generation via preference-conditioned dynamic routing,” *arXiv preprint arXiv:2502.02743*, 2025.
- [38] J. Hu, Y. Wang, S. Zhang, K. Zhou, G. Chen, Y. Hu, B. Xiao, and M. Tan, “Dynamic ensemble reasoning for llm experts,” *arXiv preprint arXiv:2412.07448*, 2024.
- [39] K. Nakaishi, Y. Nishikawa, and K. Hukushima, “Critical phase transition in large language models,” *arXiv preprint arXiv:2406.05335*, 2024.
- [40] Y. Li, J. Zhang, S. Feng, P. Yuan, X. Wang, J. Shi, Y. Zhang, C. Tan, B. Pan, Y. Hu *et al.*, “Revisiting self-consistency from dynamic distributional alignment perspective on answer aggregation,” *arXiv preprint arXiv:2502.19830*, 2025.