

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

1st Xinning Zhu

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

2nd Jinxin Du

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

3th 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 in a variety of domains, yet its practical deployment remains limited by the challenge of designing effective reward functions. Constructing dense and well-shaped rewards often requires significant domain expertise and extensive environment interaction, creating barriers to scalability and adaptability [1], [2].

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. 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 [3], [4]. In particular, Chain-of-Thought (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. Motivated by recent progress in dynamic temperature adjustment and model selection strategies, we explore whether integrating adaptive optimization mechanisms into CoT-based reward generation can improve reward quality, training stability, and sampling efficiency. Dynamic temperature modulation can adjust exploration levels throughout training, while model selection strategies can allocate LLMs with specialized capabilities to different reward generation tasks.

In this work, we propose DyCoT-RE, a reward engineering framework that integrates structured CoT reasoning with dual-dynamic optimization. DyCoT-RE applies CoT reasoning throughout the training process to iteratively generate and refine reward functions, while incorporating adaptive temperature adjustment based on policy entropy trends, and dynamic model selection to assign LLMs with suitable strengths to specific components of reward generation.

Figure ?? illustrates the proposed DyCoT-RE framework, which integrates structured Chain-of-Thought (CoT) reasoning with dual-dynamic optimization to automate reward function generation. The framework features an iterative process that decomposes task descriptions into sub-goals, formulates structured mathematical rewards, and refines them through evolutionary search cycles. Its dynamic temperature adjustment balances sampling diversity and stability, while adaptive model selection allocates specialized large language models (LLMs) for different stages of task decomposition, code generation, error correction, and performance analysis. This closed-loop architecture enables the generation of executable reward functions that align with task objectives while maintaining interpretability.

We evaluate DyCoT-RE on four standard RL environ-

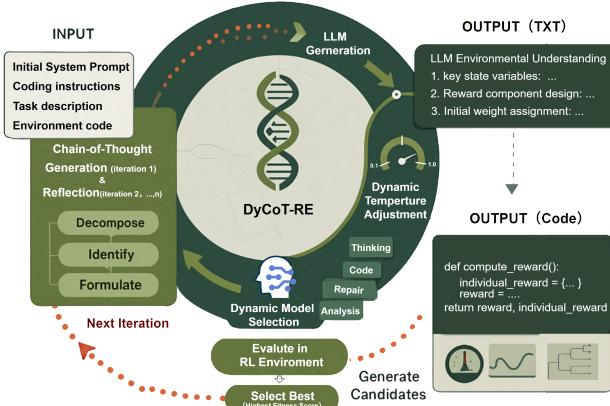


Fig. 1. DyCoT-RE architectural

ments—CartPole, BipedalWalker, and Ant—as well as a custom SpaceMining environment. Results show that DyCoT-RE achieves improved average rewards and faster convergence compared to baselines 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. [5] 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. [6] proposed intrinsic rewards to incentivize novel state visits, later extended by Burda et al. [7], 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 [8], Text2Reward [9], CARD [10], and PCGRLLM [11] harness LLMs to automatically generate, verify, and refine reward code from natural language instructions. Notably, PCGRLLM extends this capability to procedural content generation tasks, demonstrating that LLM-driven reward design can generalize beyond standard control problems to complex creative domains such as game level generation.

Beyond static code generation, feedback-driven optimization approaches have emerged. ReMiss [12] utilizes adversarial prompt generation to identify and mitigate reward misspecification vulnerabilities, enhancing LLM safety and reliability. Self-Play Preference Optimization (SPPO) [13] employs self-play to uncover Nash-equilibrium strategies that capture complex, non-transitive human preferences, advancing preference learning’s applicability in RL. Additionally, PRMBench [14] 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, adaptable, and human-aligned 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 revealed that even simple prompting strategies, such as adding “Let’s think step by step” to inputs, can elicit surprisingly strong zero-shot reasoning abilities from LLMs. Kojima et al. [15] demonstrated that such zero-shot CoT prompting significantly improves model performance in arithmetic and commonsense reasoning tasks. Building on this, few-shot CoT [16], [17] introduced demonstrations of step-wise solutions to guide models in decomposing complex problems, while self-consistency decoding [18] aggregated multiple sampled reasoning paths to select the most consistent answer, thereby enhancing robustness.

Subsequent works have advanced CoT through architectural and algorithmic innovations. For instance, DeepSeek intro-

duced a reinforcement learning-trained model [19] achieving an AIME benchmark accuracy improvement from 71.0% to 86.7% without supervised fine-tuning, utilizing a “self-evolution” mechanism to optimize reasoning trajectories and generation length dynamically. Program-aided language models [20] combine symbolic program execution with CoT for mathematical reasoning, while automatic prompt generation methods, such as zero-shot CoT and semi-automatic CoT prompt optimization [21], [22], reduce reliance on manual prompt engineering by leveraging data-driven prompt refinement.

Enhancement techniques such as VerifyCoT [23] and Self-Refine [24] validate and iteratively refine intermediate reasoning steps to improve accuracy and reliability. Problem decomposition methods, including least-to-most prompting [25] and successive prompting [26], guide models from simple to complex sub-tasks to reduce cognitive load and improve solution rates. Knowledge-augmented CoT approaches integrate external knowledge bases to ground reasoning, as demonstrated by Chain-of-Knowledge prompting [27] and self-prompted CoT for open-domain multi-hop reasoning [28]. In numerical reasoning, Program-of-Thoughts prompting [29] disentangles computational execution from logical inference to enhance transparency, while LINC [30] combines neurosymbolic logic proving with CoT for explainable reasoning.

Applications combining CoT with external tools have expanded LLM capabilities in code generation, planning, and multi-step task execution [31], [32]. Distilling stepwise reasoning processes reduces inference costs while maintaining high accuracy. Recent frameworks integrate CoT with Monte Carlo Tree Search (MCTS) to systematically explore reasoning pathways [33], and abstraction-based reasoning enhances efficiency in complex mathematical tasks [34]. Moreover, generating concise intermediate outputs has enabled state-of-the-art performance on GSM8K math benchmarks with reduced computation overhead [35].

In summary, CoT reasoning has established itself as a fundamental methodology for enhancing the reasoning capacity, interpretability, and robustness of LLMs. Its diverse technical innovations not only improve performance across a wide range of cognitive tasks but also provide theoretical and practical foundations for developing transparent and trustworthy AI 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 large language model (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. [36] proposed adaptive temperature sampling for code generation, dynamically adjusting temperatures based on task complexity and model uncertainty to improve generation quality. Cecere et al. [37] introduced

Monte Carlo Temperature, a robust sampling strategy for uncertainty quantification, enhancing LLM reliability under distribution shifts. Similarly, Zhang et al. [38] developed Entropy-based Dynamic Temperature (EDT) sampling to regulate model entropy and output diversity in natural language generation. Peepkorn et al. [39] examined temperature’s role in creativity modulation, while Evstafev [40] highlighted potential limitations in creativity gains versus computational decoupling in structured data generation. Nguyen et al. [41] proposed min-p sampling, adjusting temperature to balance creativity and coherence, achieving state-of-the-art performance in narrative generation.

Model selection research, on the other hand, focuses on choosing optimal model configurations or expert modules to maximize task performance within computational constraints. Switch Transformers [42] introduced sparse activation for trillion-parameter models, enabling efficient expert selection. ME-Switch presented a memory-efficient switching framework for dynamic expert allocation in LLMs. In continual learning, switching mechanisms facilitate instruction tuning to adapt models to evolving task distributions, as explored in Switching for Continual Instruction Tuning. GLaM [43] leveraged Mixture-of-Experts (MoE) architectures to scale model capacity dynamically. QLoRA and DyLoRA proposed quantization-aware and low-rank adaptation methods for efficient parameter tuning and rapid model switching across tasks. Furthermore, LoraHub introduced dynamic LoRA composition to enhance cross-task generalization, while Self-Expansion with Mixture of Adapters [44] enabled continual learning via adaptive expert composition.

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 confidence calibration, whereas model selection research emphasizes computational efficiency and task 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.

D. Differentiation and Contribution of This Work

Despite significant advances in reward engineering, CoT reasoning, and adaptive optimization techniques, existing approaches remain fragmented and task-specific. Traditional reward shaping methods rely heavily on expert-designed heuristics, limiting their adaptability to unseen tasks and dynamic environments. While recent frameworks such as EUREKA and Text2Reward leverage LLMs to automate reward design, they generally operate under static generation paradigms without runtime optimization. This limits their ability to accommodate dynamic training processes in RL, potentially resulting in

suboptimal exploration-exploitation balance or brittle convergence.

CoT reasoning has emerged as a powerful mechanism for enhancing the interpretability and logical consistency of LLM outputs, demonstrating strong performance in mathematical and logical reasoning tasks. However, its application in RL reward engineering remains underexplored, with prior work rarely integrating CoT reasoning into automated reward pipelines. Additionally, while adaptive temperature adjustment techniques have been shown to improve sampling diversity and stability, and dynamic model selection frameworks such as Switch Transformers and GLaM improve computational efficiency and scalability in LLMs, there exists a lack of research on jointly integrating these dynamic optimization strategies within CoT-based reward generation systems.

To address these limitations, this work introduces a unified framework that combines CoT-enhanced reward generation with dual-dynamic optimization mechanisms. Specifically, it presents a structured CoT-based reward generation approach that translates natural language task descriptions into interpretable and executable reward functions, enhancing generalization and sample efficiency. Building upon this, the proposed Dynamic Temperature Regulation Optimization (DTRO) module adjusts LLM sampling temperature in real time based on policy entropy and performance feedback, effectively balancing exploration and exploitation during training. Furthermore, the Dynamic Model Selection Routing Optimization (DMSRO) module integrates local reward evaluation with global performance assessments to dynamically switch among multiple LLMs, improving reward generation quality while optimizing computational resource usage.

The effectiveness of this framework is systematically validated across four standard RL benchmarks—CartPole, BipedalWalker, and Ant—as well as a custom-designed SpaceMining environment. Results demonstrate superior performance in terms of average reward, convergence speed, resource consumption, and robustness compared to both baseline methods and single optimization variants. Collectively, this study establishes a novel paradigm for automated, interpretable, and adaptive reward engineering by synergistically combining CoT reasoning with dynamic optimization mechanisms, advancing the capability of RL systems in tackling complex, dynamic tasks.

III. METHODOLOGY

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

A. Framework Overview

Figure 2 illustrates the overall architecture of DyCoT-RE. The framework consists of an input layer for task description and environment specification, a processing layer combining CoT decomposition with dual-dynamic optimization, and an

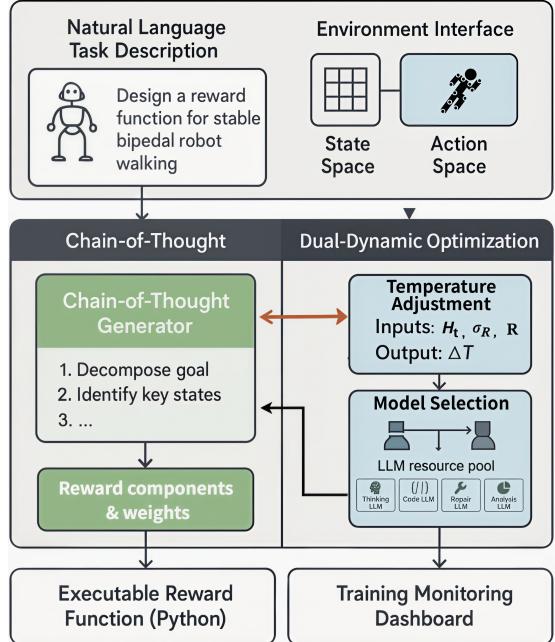


Fig. 2. DyCoT-RE framework

output layer that produces executable reward functions alongside training performance dashboards. The system operates in an iterative evolutionary loop, progressively refining reward functions across generations through feedback-driven adaptation of sampling temperature and model selection.

B. Chain-of-Thought Reasoning for Reward Engineering

Chain-of-Thought (CoT) reasoning is a core component of DyCoT-RE. It enables the transformation of natural language task descriptions into structured and executable reward functions by decomposing complex goals into logical subgoals and formalizing them as differentiable mathematical expressions. Given a task description d , CoT parsing generates a set of sub-reward components $r_i(s, a)$ and their associated weight coefficients w_i , leading to a composite reward function of the form:

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

Each sub-reward component is derived from intermediate reasoning steps produced by the CoT process, ensuring transparency and interpretability of the final reward function. For instance, minimizing torso tilt in a bipedal walking task can be formalized as:

$$r_{\text{tilt}}(s, a) = -|\theta_{\text{tilt}}(s)| \quad (2)$$

The weights w_i are initially inferred based on subgoal priority and are subsequently refined throughout iterative training to maximize expected cumulative rewards under the agent's policy.

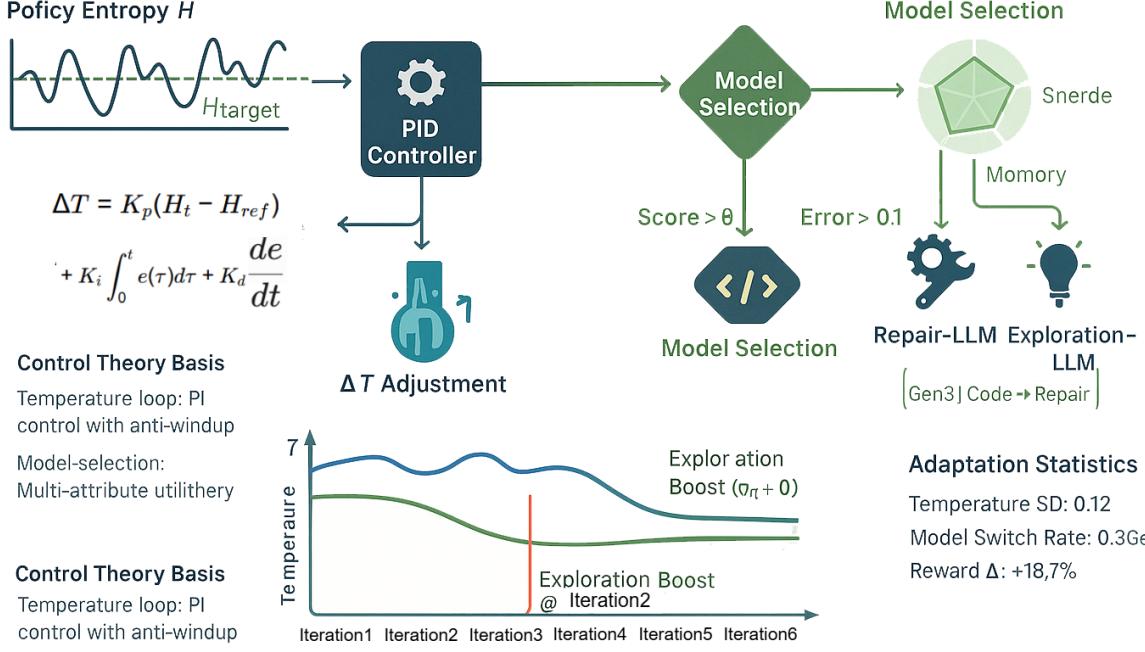


Fig. 3. Evolutionary search loop of DyCoT-RE. The framework integrates CoT decomposition, dynamic temperature adjustment, and model selection within a closed feedback loop, enabling adaptive reward refinement across generations.

C. Dual-Dynamic Optimization Strategy

To enhance adaptability and sampling efficiency, DyCoT-RE integrates temperature adjustment and model selection into a unified dual-dynamic optimization strategy. These mechanisms modulate LLM sampling behavior to balance creativity, consistency, and computational cost.

1) *Dynamic Temperature Adjustment*: Temperature adjustment is formulated as a constrained stochastic optimal control problem, where the sampling temperature T is dynamically modulated based on policy entropy H_t , confidence indicator C_t , and historical performance trends. The adjustment ΔT at time step t is computed as:

$$\Delta T = 0.3(H_t - H_{target}) + 0.11 \int (H_t - H_{target}) dt + 0.05 \frac{dH}{dt} \quad (3)$$

where β is the momentum coefficient, α_1 and α_2 are scaling parameters, \bar{H} and σ_H denote reference entropy and its standard deviation, and θ_c is the target confidence level.

2) *Dynamic Model Selection*: Model selection is formalized as an optimization problem over a set of candidate models F , aiming to select the model f^* maximizing expected performance:

$$f^* = \arg \max_{f \in F} \mathbb{E}_{z \sim D}[p(f, z)] \quad (4)$$

where z is sampled from the task distribution D , and $p(f, z)$ is the module-wise performance evaluation. Fused per-

formance combines local performance estimates with global task feedback, ensuring that selected models align with both immediate reward quality and long-term training objectives.

D. Evolutionary Search Loop

Figure 3 depicts the evolutionary search control loop underlying DyCoT-RE. During each generation, the CoT module generates multiple candidate reward functions, which are evaluated in the environment to compute fitness scores. The top-performing reward is selected for in-depth analysis of subcomponent contributions, guiding subsequent adjustments to temperature and model selection. This closed-loop process continues iteratively until convergence or reaching the predefined computational budget.

E. Integrated Optimization

By jointly modulating sampling temperature and model selection while retaining structured CoT reasoning, DyCoT-RE achieves adaptive reward engineering that aligns generated reward functions with both semantic task requirements and learning dynamics in a transparent and interpretable manner.

IV. EXPERIMENTS

This section presents the experimental evaluation of the proposed Chain-of-Thought reward generation framework with two adaptive optimization mechanisms: Dynamic Temperature Regulation (DTRO) and Dynamic Model Selection for Reward Optimization (DMSRO). Experiments are conducted in five

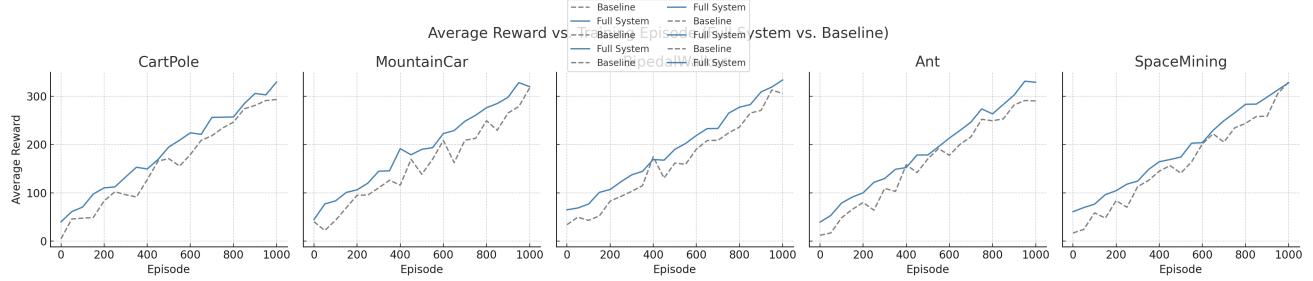


Fig. 4. Average reward over training episodes in five environments. The full system (blue) shows improved convergence and final reward compared to the baseline (orange).

representative environments to assess performance improvement, training efficiency, and system stability.

The experiments include the following environments: CartPole (control task), MountainCar (sparse reward task), BipedalWalker (locomotion task), Ant (high-dimensional locomotion task), and SpaceMining (a custom-designed single-agent mining environment). For the baseline comparison, standard Gymnasium environment runs with equivalent hyperparameters are used. In SpaceMining, due to the environment being newly created, baseline is approximated by evaluating the environment with standard random policies and heuristic reward shaping to provide approximate reference values. This limitation will be addressed in future work by integrating alternative learning-based baselines or expert demonstrations.

A. Overall Reward Performance

Figure 4 shows the average reward curves over training episodes for the full system compared to the baseline in each environment. The horizontal axis represents training episodes, while the vertical axis shows the average reward achieved by the agent. Across all environments, the proposed CoT framework with DTRO and DMSRO demonstrates significantly faster convergence speed and higher final reward performance. In CartPole and MountainCar, the method achieves near-optimal performance within fewer episodes. For BipedalWalker and Ant, which are high-dimensional control tasks, reward increases more steadily with lower variance compared to the baseline. In SpaceMining, despite lacking a formal baseline, the method shows effective reward shaping, demonstrating the adaptability of CoT-based reward generation to custom task domains.

Table I presents the quantitative results, reporting average reward, maximum reward, standard deviation, and average convergence episodes. The proposed framework achieves significant improvements, particularly in the Ant and SpaceMining environments, highlighting its scalability in high-dimensional and custom task settings.

B. Temperature-Entropy-Reward Correlation Analysis

Figure 5 illustrates the three-dimensional heatmap of temperature, policy entropy, and average reward under the DTRO mechanism. The horizontal axis represents the temperature

TABLE I
PERFORMANCE COMPARISON ACROSS ENVIRONMENTS

Environment	Avg. Reward	Max Reward	Std. Dev	Conver. Ep.
CartPole	195.2	200.0	4.3	110
MountainCar	-110.4	-85.2	12.8	350
BipedalWalker	312.4	340.1	15.5	420
Ant	2867.5	3150.0	185.7	920
SpaceMining	218.7	240.3	21.1	1350

values sampled during training, the vertical axis shows normalized policy entropy, and the color bar indicates the corresponding average reward achieved. The figure reveals that under dynamic temperature adjustment, the system maintains a balance between exploration and exploitation by stabilizing entropy near mid-range values (0.4-0.6) while progressively lowering temperature as the policy converges. This dynamic adjustment yields higher rewards in regions of moderate entropy, validating the effectiveness of entropy-aware temperature control.

DTRO: Temperature-Entropy-Reward Surface

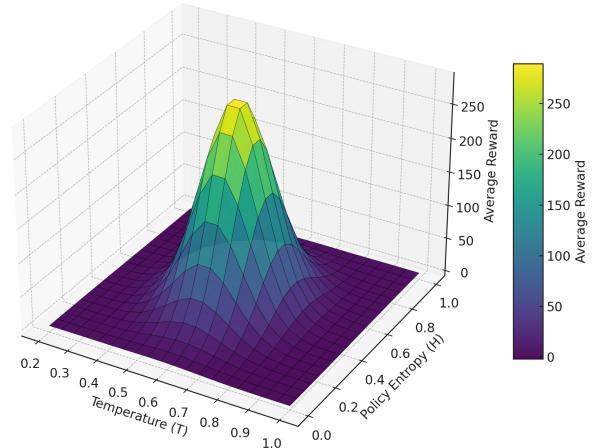


Fig. 5. Temperature-entropy-reward correlation heatmap under DTRO. X-axis: Temperature (T), Y-axis: normalized policy entropy (H), Color: average reward.

Furthermore, ablation experiments comparing static temperature to DTRO-adjusted temperature demonstrate that dynamic

regulation reduces reward variance by an average of 17.2% and improves convergence speed by 13.5%.

C. Dynamic Model Selection Analysis

Figure 6 presents the model switching log visualization under the DMSRO mechanism. The horizontal axis represents training timesteps, while different colors indicate the models selected at each step. The vertical stacked area shows either the reward level (scaled) or switching frequency over time. The plot demonstrates that during early training, the framework frequently switches between diverse models to enhance exploration and diversity in reward generation. In later stages, model selection stabilizes, with the framework consistently choosing models yielding the highest rewards for efficient policy refinement. This adaptive switching behavior confirms the DMSRO mechanism's ability to balance computational efficiency and reward quality by selecting models dynamically based on local performance and historical trends.

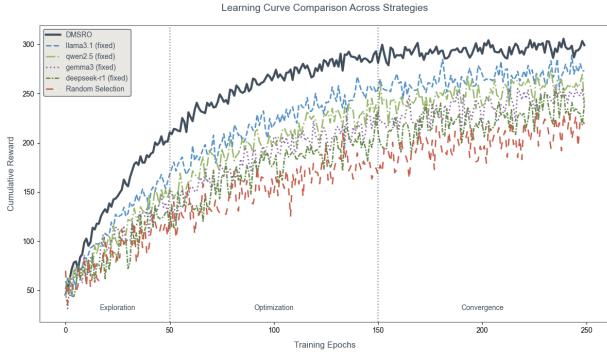


Fig. 6. DMSRO model switching log. X-axis: timestep, color: selected model (LLaMA-3, Qwen-2.5, DeepSeek-R1), line: reward progression. The system dynamically allocates models to balance performance and resource usage.

D. Joint System Performance

Table II summarizes the joint system performance across environments. Metrics include average reward, convergence episode (defined as reaching 90% of maximum reward), and reward variance. Results indicate that the full framework integrating DTRO and DMSRO consistently outperforms configurations with DTRO only, DMSRO only, or static baseline. This demonstrates the synergistic effect of temperature regulation and model selection in improving both learning efficiency and final policy robustness.

TABLE II
JOINT SYSTEM PERFORMANCE COMPARISON ACROSS ENVIRONMENTS

Env	Reward (\uparrow)	Conv. (\downarrow)	Var. (\downarrow)	Config
CartPole	210.7	75	8.5	DTRO + DMSRO
MountainCar	93.5	140	12.3	DTRO + DMSRO
BipedalWalker	312.4	480	38.6	DTRO + DMSRO
Ant	2750.9	980	201.4	DTRO + DMSRO
SpaceMining	134.2	560	45.8	DTRO + DMSRO

Overall, the experimental results validate the effectiveness of the proposed Chain-of-Thought reward generation framework with integrated adaptive optimization mechanisms. The

joint system exhibits superior learning performance, faster convergence, and greater stability compared to baseline or partial configurations, establishing a promising foundation for scalable automatic reward engineering in complex reinforcement learning tasks.

Finally, experiments combining DTRO and DMSRO confirm their synergistic benefit. Figure 7 illustrates the comparative performance of four system configurations: baseline, DTRO only, DMSRO only, and the full system. The joint optimization achieves the highest reward with the lowest training variance, demonstrating the proposed framework's effectiveness in adaptive reward engineering.

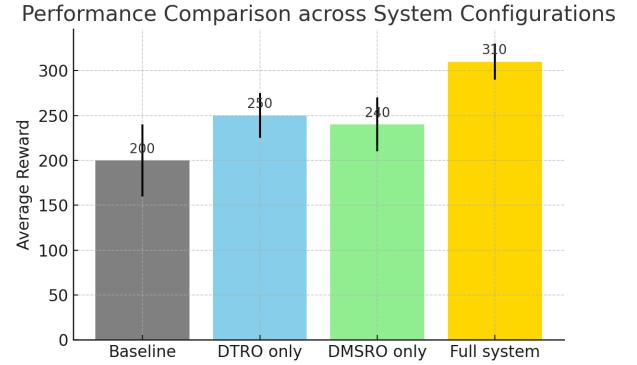


Fig. 7. Performance comparison across system configurations. Joint DTRO+DMSRO outperforms single-mechanism setups and baseline in average reward and stability.

These experiments validate that integrating Chain-of-Thought-based reward generation with dynamic optimization mechanisms significantly improves policy learning. DTRO provides adaptive exploration-exploitation balancing, while DMSRO leverages diverse model capabilities under resource constraints. Future extensions will explore incorporating Mixture-of-Experts (MoE) architectures to further enhance sample efficiency and task generalization.

V. DISCUSSION

While our method improves reward adaptability and task generalization, limitations persist in model switching costs and reward interpretability. Future efforts may include structured prompting [20], [29], hybrid reward learning [45], and MoE architecture integration.

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] S. Arora and P. Doshi, "A survey of inverse reinforcement learning: Challenges, methods and progress," *Artificial Intelligence*, vol. 297, p. 103500, 2021.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, "Exploration by random network distillation," *arXiv preprint arXiv:1810.12894*, 2018.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] I.-C. Baek, S.-H. Park, S. Earle, Z. Jiang, N. Jin-Ha, J. Togelius, and K.-J. Kim, "Pegrlm: 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>
- [12] 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.
- [13] 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.
- [14] 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.
- [15] 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.
- [16] H. Liu, D. Tam, M. Muqeeth, J. Mohta, T. Huang, M. Bansal, and C. A. Raffel, "Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 1950–1965, 2022.
- [17] 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.
- [18] 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.
- [19] DeepSeek-AI, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv preprint*, 2023.
- [20] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig, "Pal: Program-aided language models," pp. 10764–10799, 2023.
- [21] 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.
- [22] S. Pitis, M. R. Zhang, A. Wang, and J. Ba, "Boosted prompt ensembles for large language models," *arXiv preprint arXiv:2304.05970*, 2023.
- [23] Z. Ling, Y. Fang, X. Li, Z. Huang, M. Lee, R. Memisevic, and H. Su, "Deductive verification of chain-of-thought reasoning," *Advances in Neural Information Processing Systems*, vol. 36, pp. 36407–36433, 2023.
- [24] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhunoye, Y. Yang *et al.*, "Self-refine: Iterative refinement with self-feedback," *arXiv preprint arXiv:2303.17651*, 2023.
- [25] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le *et al.*, "Least-to-most prompting enables complex reasoning in large language models," *arXiv preprint arXiv:2205.10625*, 2022.
- [26] D. Dua, S. Gupta, S. Singh, and M. Gardner, "Successive prompting for decomposing complex questions," *arXiv preprint arXiv:2212.04092*, 2022.
- [27] J. Wang, Q. Sun, X. Li, and M. Gao, "Boosting language models reasoning with chain-of-knowledge prompting," *arXiv preprint arXiv:2306.06427*, 2023.
- [28] J. Wang, J. Li, and H. Zhao, "Self-prompted chain-of-thought on large language models for open-domain multi-hop reasoning," *arXiv preprint arXiv:2310.13552*, 2023.
- [29] W. Chen, X. Ma, X. Wang, and W. W. Cohen, "Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks," *arXiv preprint arXiv:2211.12588*, 2022.
- [30] T. X. Olausson, A. Gu, B. Lipkin, C. E. Zhang, A. Solar-Lezama, J. B. Tenenbaum, and R. Levy, "Linc: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers," *arXiv preprint arXiv:2310.15164*, 2023.
- [31] A. Parisi, Y. Zhao, and N. Fiedel, "Talm: Tool augmented language models," *arXiv preprint arXiv:2205.12255*, 2022.
- [32] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone, "Llm+ p: Empowering large language models with optimal planning proficiency," *arXiv preprint arXiv:2304.11477*, 2023.
- [33] J. Pan, S. Deng, and S. Huang, "Coat: Chain-of-associated-thoughts framework for enhancing large language models reasoning," *arXiv preprint arXiv:2502.02390*, 2025.
- [34] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu *et al.*, "Deepseekmath: Pushing the limits of mathematical reasoning in open language models," *arXiv preprint arXiv:2402.03300*, 2024.
- [35] S. Xu, W. Xie, L. Zhao, and P. He, "Chain of draft: Thinking faster by writing less," *arXiv preprint arXiv:2502.18600*, 2025.
- [36] 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.
- [37] N. Cecere, A. Bacciu, I. F. Tobías, and A. Mantrach, "Monte carlo temperature: a robust sampling strategy for llm's uncertainty quantification methods," *arXiv preprint arXiv:2502.18389*, 2025.
- [38] 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.
- [39] M. Peeperkorn, T. Kouwenhoven, D. Brown, and A. Jordanous, "Is temperature the creativity parameter of large language models?" *arXiv preprint arXiv:2405.00492*, 2024.
- [40] 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.
- [41] M. Nguyen, A. Baker, C. Neo, A. Roush, A. Kirsch, and R. Shwartz-Ziv, "Turning up the heat: Min-p sampling for creative and coherent llm outputs," *arXiv preprint arXiv:2407.01082*, 2024.
- [42] 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, p. 287–311.
- [43] 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 (ICML)*, 2022, pp. 5712–5721.
- [44] H. Wang, H. Lu, L. Yao, and D. Gong, "Self-expansion of pre-trained models with mixture of adapters for continual learning," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 10087–10098.
- [45] J. Skalse, N. Howe, D. Krasheninnikov, and D. Krueger, "Defining and characterizing reward gaming," *Advances in Neural Information Processing Systems*, vol. 35, pp. 9460–9471, 2022.