
RF-Agent: Automated Reward Function Design via Language Agent Tree Search

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Designing efficient reward functions for low-level control tasks is a challenging
2 problem. Recent research aims to reduce reliance on expert experience by using
3 Large Language Models (LLMs) with task information to generate dense reward
4 functions. These methods typically rely on training results as feedback, iteratively
5 generating new reward functions with greedy or evolutionary algorithms.
6 However, they suffer from poor utilization of historical feedback and inefficient
7 search, resulting in limited improvements in complex control tasks. To address
8 this challenge, we propose RF-Agent, a framework that treats LLMs as language
9 agents and frames reward function design as a sequential decision-making process,
10 enhancing optimization through better contextual reasoning. RF-Agent integrates
11 Monte Carlo Tree Search (MCTS) to manage the reward design and optimization
12 process, leveraging the multi-stage contextual reasoning ability of LLMs. This
13 approach better utilizes historical information and improves search efficiency to
14 identify promising reward functions. Outstanding experimental results in 17 diverse
15 low-level control tasks demonstrate the effectiveness of our method.

16 1 Introduction

17 Reward design is a crucial component in Reinforcement Learning (RL), significantly affecting
18 the performance and training efficiency of the policy, especially in low-level control tasks like
19 locomotion[39] and complex manipulation[3, 15]. Although task evaluation metrics (e.g., success
20 rate, movement speed) can be directly used as rewards, their sparse or one-dimensional nature
21 presents challenges for policy optimization. Therefore, reward shaping[33] through dense reward
22 functions is crucial for guiding policy learning more efficiently[45]. The most common approach
23 involves human experts manually crafting a dense reward function, which is interpretable but requires
24 extensive expertise and may be suboptimal[4]. To address this, Inverse RL[64, 53, 13] and Preference-
25 based RL[10, 21, 34] attempt to generate dense rewards by learning reward models from expert
26 demonstrations or preferences. However, they still rely on large amounts of expert data and lack
27 interpretability.

28 Recently, large language models (LLMs) have excelled in not only traditional NLP tasks[32] but also
29 logical reasoning and agent-based tasks[49, 57]. Given their impressive world knowledge [51] and
30 coding capabilities [55, 59], a natural thought is to leverage LLMs to generate interpretable dense
31 reward functions, replacing expert effort. Approaches like L2R[58] and Text2Reward[54] made
32 strides in this direction, achieving remarkable results. In addition, more recent methods[27, 17, 60, 61]
33 have further advanced the field by combining more refined feedback design and evolutionary concepts.
34 We view these methods as part of a sequential decision process, as shown in Fig1.a. The interaction
35 between the policy and environment forms an internal system, where the LLM, as a language agent,
36 continuously interacts with this system. The action of the LLM represents the process of generating

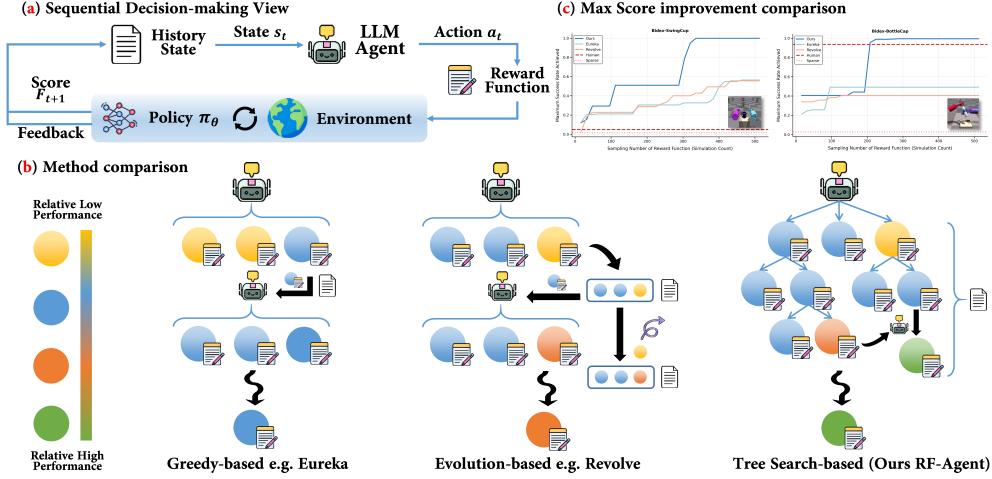


Figure 1: Quick view. (a) Reshaping LLM-based Reward Function Design Problem from sequential decision-making view. (b) Abstract pipeline comparison between different approaches. (c) Comparison between different methods on max success rates with the numbers of complete training times.

37 a reward function, and the internal system provides results after training as feedback based on the
 38 action. The action and feedback in the current decision-making process are finally integrated into the
 39 history state and utilized by LLM for future actions.

40 From the sequential decision-making perspective, current research can be categorized into two reward
 41 generation modes, as shown in Fig1.b. The first mode, exemplified by Eureka[27], follows a greedy
 42 algorithm[47] approach, where the best reward function is retained after a batch generation of reward
 43 functions and serves as the starting point for the next iteration. The second mode, represented by
 44 Revolve[17], adopts an evolutionary algorithm[11] approach, maintaining a population and evolving
 45 new reward functions by integrating individual information from the population. However, when
 46 utilizing LLMs to design reward functions for complex control tasks[9], the performance improvement
 47 of these methods remains limited, as illustrated in Fig1.c. This limitation hinders the generation of
 48 high-quality and training-efficient reward functions through LLMs, constraining the performance
 49 of this pipeline. We identify two key factors leading to limited performance improvement. First,
 50 from the search perspective, greedy or evolutionary algorithms struggle in complex decision-making
 51 problems, as they fail to effectively balance exploitation and exploration, often resulting in premature
 52 convergence to local optima or excessive exploration of possibilities. Second, regarding information
 53 usage, existing methods only retain local historical information, overlooking the potential decision
 54 paths that could transition from low-performing to high-performing reward functions (Fig1.b ours).

55 Based on the analysis above, we introduce a novel reward function design framework, RF-Agent,
 56 which treats the LLM as an agent to address the challenges outlined by enhancing their reasoning
 57 and decision-making capabilities. Specifically, we employ a tree structure to represent the entire
 58 reward design process (Fig1.b ours) and its history states. Each node in this tree corresponds to a
 59 distinct decision strategy and its associated reward function, effectively leveraging the multi-stage
 60 contextual reasoning capabilities of LLMs. To further improve the framework, we incorporate the
 61 Monte Carlo Tree Search (MCTS) algorithm[5] to balance exploration and exploitation within the
 62 decision space. Additionally, during the action expansion phase, we integrate a set of heuristic
 63 action types to guide the LLM in generating reward functions with varying inclinations. This enables
 64 RF-Agent to appropriately incorporate global decision information from other paths, ensuring the
 65 full utilization and refinement of promising reward functions. Finally, environmental feedback and
 66 self-verify metrics contribute to an overall evaluation, which forms a value prior for node selection,
 67 thereby enhancing the search optimization capabilities of the entire RF-Agent framework.

68 We evaluate the proposed RF-Agent across 17 tasks in IsaacGym[30] and Bi-DexHands[9], covering
 69 task types such as locomotion and robot arm/hand manipulation. The results indicate that our approach
 70 outperforms the current state-of-the-art LLM-based reward function design methods, producing high-
 71 performance and training-efficient reward functions. Such achievement demonstrates the potential of
 72 leveraging language agents in combination with reasoning and search frameworks for reward function
 73 design. Ablation studies further validate the effectiveness of our architecture.

74 **2 Related Works**

75 **Reward Design with LLMs.** Reward engineering has long been a challenge in reinforcement
76 learning[44, 25], with the common approach being manual creation of reward functions by experts[22,
77 4]. Inverse RL[1, 18, 64, 53, 13] and Preference-based RL[10, 21, 34] generate rewards through
78 learned models, but still require expert demonstrations and lack interpretability. Recent advancements
79 leverage LLMs, which exhibit strong world knowledge and logical reasoning[57, 14], in tasks like
80 code generation[8, 38] and robotic control planning[26, 43, 48]. Some works have extended this
81 to generating rewards for RL. For instance, [23] uses LLMs to generate binary rewards based on
82 user intent in negotiation games, and [12] calculates rewards based on the similarity between state
83 transitions and goals generated by LLMs. However, these methods suffer from excessive LLM calls
84 and inconsistencies. In contrast, L2R[58] and T2R[54] generate dense, interpretable reward functions,
85 while T2R directly offers plug-and-play reward function Python code, introducing a promising
86 paradigm for reward engineering. Recent works like Eureka[27] use greedy iterative approaches to
87 optimize reward functions, while Revolve[17] incorporates evolutionary algorithms to evolve new
88 rewards. However, both methods still face challenges with inefficient use of historical feedback and
89 low search efficiency. Our RF-Agent addresses these issues by introducing Monte Carlo Tree Search
90 to manage reward function design and optimization process from a decision-making perspective. It
91 fully utilizes LLMs’ contextual reasoning and historical feedback, thereby significantly enhancing
92 the performance of the designed reward function.

93 **LLM-based Agent for Decision Making.** The world knowledge and logical reasoning abilities
94 of LLMs enable them to serve as agents for solving decision-making problems[35, 48, 6]. For
95 example, LLMs have been used as high-level controllers in robotic planning [2], interacted with web
96 browsers to fulfill user needs [31], and tackled multi-step decisions in text-based adventure games
97 [42]. As tasks become more complex, methods such as Twosome[46] attempt to fine-tune LLM-based
98 agents to interact directly with the environment, but this often incurs high training costs. In contrast,
99 training-free approaches that utilize multi-step reasoning to leverage LLMs’ contextual learning have
100 gained popularity[52, 50, 57, 62, 37]. This “Test-time Compute” paradigm [24] enhances reasoning to
101 influence decision performance, with Chain-of-Thought (CoT)[52] being the most prominent method.
102 ReAct[57] incorporates CoT into decision-making, while self-refine[28] and Reflexion[41] add self-
103 feedback mechanisms. However, these linear methods neglect alternative decision paths. ToT[56]
104 breaks the decision process down using tree search, and methods like LATS[63] and RAP[16] further
105 introduce Monte Carlo Tree Search (MCTS) to enhance reasoning during decision-making. Inspired
106 by these approaches, we frame reward function design as a decision-making process, enhancing the
107 logical reasoning of LLM to better analyze historical feedback and generate high-quality reward
108 functions.

109 **3 Problem Setting**

110 We first introduce the **Reward Design Problem** (RDP)[44], given by the tuple $P = \{M, \mathcal{R}, \pi, F\}$,
111 where $M = (S, A, T)$ is the environment world model including state space S , action space A , and
112 transition function T . \mathcal{R} indicates the space of all reward functions where each $R \in \mathcal{R}$ is a reward
113 function that maps a state-action pair to a scalar $R : S \times A \rightarrow \mathbb{R}$. $\pi : S \rightarrow \Delta A$ indicates the policy
114 and $F(\cdot) : \Pi \rightarrow \mathbb{R}$ is the evaluation metric that measures the real performance of π on a given task.
115 Given a reward function R , the tuple (M, R) forms a *Markov Decision Process*(MDP)[36]. We
116 additionally denote $\mathcal{A}_M(\cdot) : \mathcal{R} \rightarrow \Pi$ as a learning algorithm (e.g., PPO[40]) that outputs the policy π
117 under the given MDP (M, R) . The goal of the RDP is to output a reward function $R \in \mathcal{R}$ such that
118 the policy $\pi := \mathcal{A}_M(R)$ under the MDP (M, R) achieves the highest evaluation score $F(\pi)$.

119 In our **Reward Function Design Problem** (RFDP) setting, each $R \in \mathcal{R}$ is specified as Python code
120 and given by a reward designer $G : \mathcal{L} \rightarrow \mathcal{R}$ according to the input string L as instruction. Here we
121 choose a pre-trained language model $p_\theta(\cdot)$ as the reward function designer G . After the policy π
122 completes the training of the fix number of steps under the given reward function R and learning
123 algorithm \mathcal{A}_M , we will receive feedback from the environment about this learning process. In
124 addition to the evaluation score F , a string $l_{feedback}$ as language feedback will indicate the execution
125 situation. Finally, given a string l_{task} that specified the task, the objective of the RFDP is to output a
126 reward function code $R \sim p_\theta(l)$ such that $F(\mathcal{A}_M(R))$ is maximized.

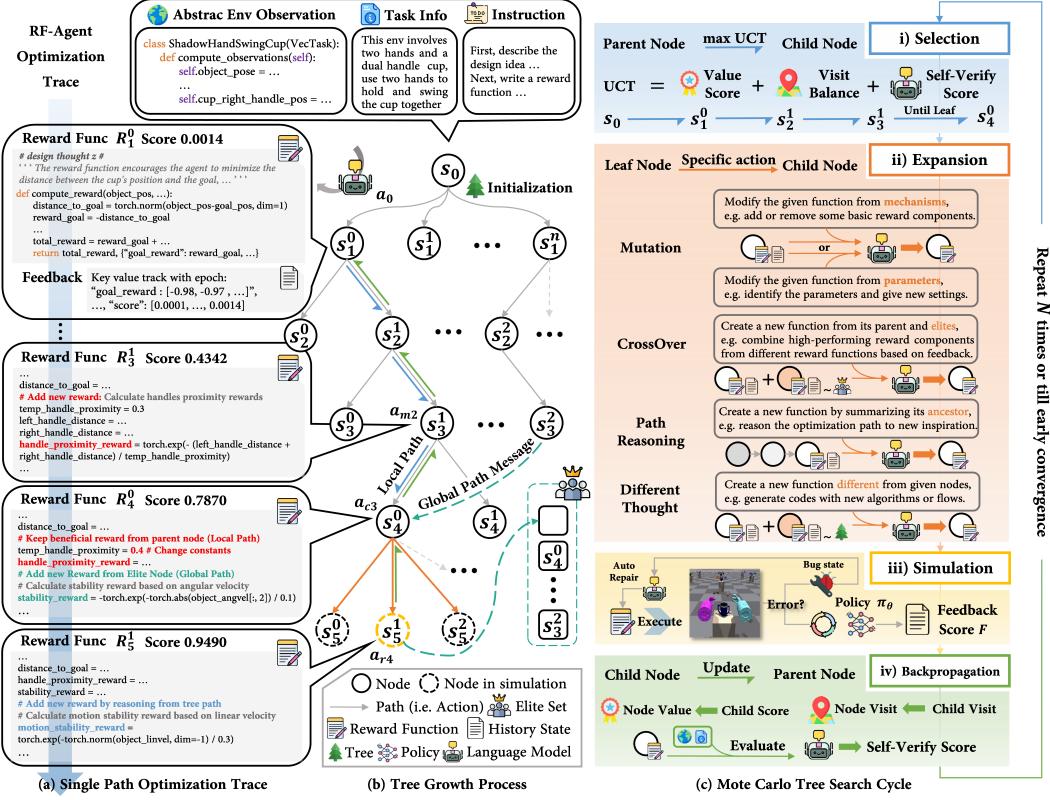


Figure 2: Illustration of our RF-Agent with (a) an exemplary reward function optimization path from the tree, (b) total tree growth process, and (c) iteration based on MCTS. Please refer to the Appendix B for prompts used in RF-Agent, mainly including initialization, expansion, and other process.

127 4 Method

128 In this section, we first explore the challenges associated with solving the LLM-based RFDP problem,
129 with a particular emphasis on the in-context learning capabilities of LLMs. We then introduce RF-
130 Agent, a framework that enhances the multi-stage contextual reasoning of LLMs by integrating Monte
131 Carlo Tree Search (MCTS), thereby improving the efficiency of reward function search through better
132 utilization of historical feedback and evaluation metrics. Finally, we show an exemplary reward
133 function optimization process to illustrate the effectiveness of RF-Agent.

134 4.1 Rethinking LLM-Based Reward Function Design Problem

135 From a decision-making perspective, each reward function R generated by an LLM $p_\theta(l)$ can be
136 viewed as an action a , with the corresponding input instruction l as the state s . This presents a key
137 challenge in RFDP: the state and action space defined on language is vast, while the evaluation
138 metrics F and feedback $l_{feedback}$ provided by the environment after training are sparse and text-based.
139 As a result, using feedback to directly optimize $p_\theta(\cdot)$ via gradients is impractical and costly.

140 The in-context learning ability of LLMs offers a promising alternative to solve the above problem. In
141 many text-based tasks, LLMs act as language agents, leveraging varying state information as context
142 to iteratively make decisions[57, 63]. Based on the analysis, we identify that previous LLM-based
143 RFDP methods simply treat LLM contextual learning as a basic usage, failing to integrate it effectively
144 within the RFDP framework. These methods suffer from limited historical information utilization
145 and low search efficiency to generate high-performance reward functions. To address these issues, we
146 propose the RF-Agent framework, which fully exploits multi-stage contextual reasoning of LLM to
147 better utilize evaluation metrics and textual feedback in RFDP and improve the efficiency of reward
148 function search.

149 **4.2 RF-Agent for Reward Function Design**

150 **Overview.** Our RF-Agent is illustrated in Fig2. We structure the optimization to solve the LLM-based
 151 RFDP as a Monte Carlo Tree Search process, with the root node n_{root} serving as a virtual node. The
 152 remaining nodes each represent a reward function R with its training feedback. Based on classical
 153 MCTS, we first generate an initial set of child nodes, after which the RFDP solution process consists
 154 of four stages: selection, expansion, simulation, and backpropagation. These stages iterate until the
 155 maximum iteration limit N is reached or early convergence. In the selection phase, we select the most
 156 promising node based on an improved Upper Confidence Bound for Trees (UCT), which incorporates
 157 evaluation metrics, self-verify metrics, and exploration count. During the expansion phase, RF-Agent
 158 utilizes specific actions to guide the LLM in generating reward functions for child nodes, ensuring
 159 effective exploration of the reward function space. Finally, in the simulation and backpropagation
 160 phase, RF-Agent trains policies with generated reward functions and receives feedback to update the
 161 node states for next selection. See Appendix G for the algorithm of RF-Agent.

162 **Initialization.** During the initialization phase, in addition to the brief task information l_{task} , RF-
 163 Agent is only granted access to the observational part of the environment code, l_{obs} , consistent
 164 with previous approaches[27, 17]. We then provide both inputs to the LLM, leveraging its world
 165 knowledge and zero-shot instruction-following capabilities to generate the reward function:

$$R \sim p_{\theta}(x, z), z \sim p_{\theta}(x), \text{ where } x = [l_{task}, l_{obs}, l_{action}] \quad (1)$$

166 Here, l_{action} refers to the instructions, and z represents the design thought in the reasoning process.
 167 As shown in Fig.2.a, we prompt the LLM to first generate the design thought and then create the
 168 specific reward function based on it, leveraging the advantages of logical reasoning[52].

169 The generated reward functions are then integrated with the environment to train the policy under
 170 the learning algorithm \mathcal{A}_M . The process returns an evaluation metric F (e.g., task success rate) and
 171 feedback $l_{feedback}$, which tracks the changes of key components in R during training as shown in
 172 App. B.2. Then, the construction from the root node n_{root} to the first layer of child nodes is complete,
 173 with each node state represented as $s = [z, R, F, l_{feedback}]$. The tree structure stores all of the node
 174 states throughout the development, enabling their full utilization in the selection and expansion stages.

175 **Selection to the Potential Node.** In the selection stage, RF-Agent balances exploration and exploita-
 176 tion to identify promising nodes as shown in Fig2.c.i and blue arrow in b, improving search efficiency
 177 for promising reward functions. We complete this stage by sequentially selecting child nodes with
 178 the highest UCT values until a leaf node. In the RFDP setting, each node maintains a value $Q(s)$
 179 reflecting the score of its reward function and those of its descendants, along with a visit count $N(s)$.

180 Additionally, RF-Agent introduces a self-verify score to improve selection. In complex control tasks,
 181 early-generated reward functions often have low scores F , even nearly zero, resulting in the value
 182 Q constructed based on F not being able to effectively reflect the potential of the node. Inspired by
 183 the inherent evaluation capabilities of LLMs[41], we prompt the LLM to first analyze how an expert
 184 completes the task, then output a score representing the likelihood of obtaining an expert-level policy
 185 by training under the current R as shown in App. B.4. The final UCT calculation is as follows:

$$\text{UCT}(s_{child}) = \frac{Q(s_{child}) - Q_{min}}{Q_{max} - Q_{min}} + \lambda \cdot \left(\sqrt{\frac{2 \cdot \log(N(s_{parent}) + 1)}{N(s_{child})}} + \sigma(v_{self}(s_{child})) \right) \quad (2)$$

186 Here, we normalize the Q -value to enhance the robustness of the UCT calculation in different tasks.
 187 The parameter λ controls the weight of the exploration term and linearly decays during the search to
 188 promote early exploration while ensuring later convergence. v_{self} represents the self-verify score
 189 scaling with softmax σ and also decays with λ .

190 **Expansion with LLM-based actions.** In the expansion stage, RF-Agent incorporates historical
 191 information to guide the LLM in generating new reward functions similar to Eq.1, but with different
 192 l_{action} . Previous methods[27] typically regenerate reward functions based on a single prompt and
 193 limited historical data, thus cannot fully leverage the LLM in-context learning abilities and limit
 194 effective exploration of the reward function space. When humans tackle problems, a complete
 195 decision-making process usually involves various types of action, with different actions employed
 196 under different states. Inspired by this, RF-Agent enhances exploration of the reward function space
 197 by defining different action types to generate the reward functions with utilization of historical
 198 information from the whole tree, as shown in Fig2.c.ii and orange arrow in b. See App. B.3 and I.1
 199 for detailed prompts and effects. The action types including $[a_{m1}, a_{m2}, a_{c3}, a_{r4}, a_{d5}]$ are as follows:

- 200 \diamond *Mutation* a_{m1} & a_{m2} : $l_{action} = l_{mutation}(s_{parent})$. We introduce the mutation to represent local
 201 modifications to the existing reward function. Given a parent node, the mutation action guides the
 202 LLM to analyze the state of the parent node to optimize its reward function. Specifically, a_{m1} guides
 203 the LLM to modify the structure of the reward function, such as adding or removing a component,
 204 while a_{m2} focuses on adjusting the parameter weights within the reward function.
- 205 \diamond *Crossover* a_{c3} : $l_{action} = l_{crossover}(s_{parent}, s_{node \sim \{elites\}})$. To enhance the utilization of global
 206 historical information and accelerate the search process, we introduce a crossover operation to extract
 207 information from high-performance nodes. RF-Agent maintains an elite set, dynamically storing
 208 high-performance nodes from the tree based on evaluation scores F . We sample k nodes from the
 209 elite set, weighted by their scores, and input them along with the parent node state to the LLM. Then
 210 RF-Agent guides the LLM to analyze the feedback and extract useful reward components from these
 211 nodes, combining them to form a new reward function.
- 212 \diamond *Path Reasoning* a_{r4} : $l_{action} = l_{reason}(s_{parent}, s_{ancestor})$. Each path from the root node to a leaf
 213 node represents a unique reward function optimization trace. Similar to how humans recall past
 214 experiences when solving new problems, RF-Agent leverages the knowledge in the reward function
 215 optimization trace by reasoning. By truncating a k -length tree path from the current node to the root,
 216 RF-Agent guides the LLM to reason the reward function optimization process, identifying the design
 217 strengths and further generating a new reward function.
- 218 \diamond *Different Thought* a_{d5} : $l_{action} = l_{differ}(s_{parent}, s_{node \sim \{tree\}})$. To prevent premature convergence
 219 during the search and enhance exploration of the reward function space, we introduce a_{d5} to
 220 generate thoughts distinct from existing reward functions in the tree. RF-Agent randomly selects
 221 nodes from different paths, combines them with the parent node, and guide the LLM to generate a
 222 reward function structurally different from those in the selected nodes.

223 **Simulation with Reward Function.** In the simulation stage, we train the policy using the reward
 224 function R of expanded nodes, as shown in Fig. 2.c.iii. Generally, R can be executed at once, then
 225 a complete node $s = [z, R, F, l_{feedback}]$ will be established. In complex environments, however,
 226 early LLM-generated reward functions may encounter execution errors. In such cases, we extract the
 227 traceback error and prompt the LLM to adjust R . This adjustment may cause the reward function \hat{R}
 228 to diverge from the original design thought z , and even without adjustments, LLM hallucinations
 229 can misalign R and z [19]. Since expansion actions are partly based on z , its accuracy is crucial. To
 230 address this, RF-Agent introduces a simple yet effective thought alignment process: during expansion,
 231 we generate z before R ; in simulation, once the reward function compiles correctly, we regenerate a
 232 more complete design idea z_{new} based on R or \hat{R} . See Appendix B.4 and I.2 for details.

233 **Backup with Evaluation.** The purpose of the backpropagation stage is to update the value $Q(s)$ and
 234 visit count $N(s)$, with generating a self-verify score v_{self} , as shown in Fig2.c.iv and green arrow in
 235 b. The process and motivation for generating the self-verify score are detailed in the expansion part,
 236 while the update process for $Q(s)$ and $N(s)$ with update rate η is as follows:

$$Q(s) = (1 - \eta) \cdot Q(s) + \eta \cdot \max_{s' \in \text{Child}(s)} Q(s'); N(s) = \sum_{s' \in \text{Child}(s)} N(s') \quad (3)$$

237 4.3 Optimization of Reward Function in RF-Agent

238 As shown in Fig2.a, we illustrate an exemplary reward function optimization path for Swing Cup task.
 239 During early exploration, RF-Agent tends to expand extensively to try different reward functions
 240 such as using mutation actions to add the handle distance reward in R_3^1 . Later, RF-Agent focuses
 241 on exploiting a few high-performance paths, such as continuing with R_3^1 by crossover of elite nodes
 242 to generate R_4^0 , which modifies parameters and incorporates angular velocity stabilization rewards.
 243 Further optimization through path reasoning led to R_5^1 , which adds a linear velocity stabilization
 244 reward and achieves high performance. See Appendix F.3 for more optimization examples.

245 5 Experiments

246 5.1 Experiment Environments

247 We test RF-Agent in two low-level control environments: IsaacGym[30] and Bi-DexHands[9],
 248 encompassing 8 control agents and 17 diverse tasks. We first evaluates on 7 representative tasks

Table 1: Quantitative evaluations on IsaacGym. **Bold** indicates the best results within the same group. Avg norm score represents the average performance after calculating human normalized scores (*Method – Sparse*) / (*Human – Sparse*) for each task.

task	locomotion				manipulation			Avg norm score
	Ant	Anymal	Humanoid	Quadcopter	AllegroHand	FrankaCabinet	ShadowHand	
Sparse	0.14 \pm 0.05	–2.05 \pm 0.24	3.01 \pm 1.99	–1.35 \pm 0.04	0.03 \pm 0.00	0.04 \pm 0.08	0.04 \pm 0.00	0
Human	6.75 \pm 0.30	–0.03 \pm 0.00	5.89 \pm 0.81	–0.08 \pm 0.03	11.41 \pm 1.62	0.11 \pm 0.07	8.56 \pm 1.23	1
LLM-based Reward Function Design with GPT-4o-mini								
Eureka	4.87 \pm 1.34	–0.67 \pm 0.04	3.29 \pm 0.19	–0.07 \pm 0.02	13.87 \pm 3.04	0.01 \pm 0.00	10.81 \pm 1.06	0.63
Revolve	5.16 \pm 1.37	–0.80 \pm 0.09	3.12 \pm 0.33	–0.08 \pm 0.02	18.40 \pm 0.97	0.19 \pm 0.17	10.61 \pm 1.38	0.67
Ours	6.22\pm0.90	–0.01\pm0.00	5.91\pm1.21	–0.03\pm0.00	22.52\pm1.83	0.35\pm0.15	13.13\pm0.82	1.70
LLM-based Reward Function Design with GPT-4o								
Eureka	5.77 \pm 0.85	–0.01 \pm 0.00	5.44 \pm 0.88	–0.03 \pm 0.02	22.55 \pm 0.99	0.52 \pm 0.21	12.27 \pm 0.91	2.00
Revolve	5.90 \pm 0.82	–0.01 \pm 0.00	4.87 \pm 0.23	–0.04 \pm 0.02	23.11 \pm 1.10	0.57 \pm 0.50	9.02 \pm 1.00	2.03
Ours	7.10\pm0.11	–0.01\pm0.00	6.37\pm0.32	–0.03\pm0.01	24.04\pm1.41	0.79\pm0.18	14.09\pm0.40	2.68

249 across locomotion and robot arm/hand manipulation from IsaacGym. To further test our method on
250 more complex tasks, we selected 10 tasks from the Bi-DexHands, involving dual-arm manipulation
251 with more intricate actions such as door closing, cup rotation, and bottle cap twisting. We divided
252 these 10 tasks into two groups, expert-easy and expert-hard, based on human reward function success
253 rates, enabling a comprehensive comparison of different LLM-based RFDP methods and human
254 performance. Task details are provided in the Appendix A.

255 5.2 Baselines

256 **Human.** These use dense reward functions from the benchmark, written by reinforcement learning
257 researchers who designed the tasks, representing expert-level reward engineering techniques.

258 **Sparse.** These are equivalent to using the evaluation score F , which directly reflects task completion
259 quality, as the reward. This reward is typically sparse or single-dimensional.

260 **Eureka**[27] automatically generates dense reward functions by combining task information and
261 environmental observations. More importantly, it generates rewards in batches and retains the best
262 one based on feedback, optimizing the reward function through a greedy iterative approach.

263 **Revolve**[17] maintains a population and applies evolutionary operations to modify the reward
264 functions with human-in-the-loop feedback. In our fully automated reward function design scenario,
265 we implemented the Revolve auto version, where feedback is also provided by the environment.

266 5.3 Training Setup

267 **Policy Training.** Following Eureka, all tasks are implemented with a well-tuned PPO[29] and default
268 hyperparameters provided by the benchmark. We tested the final reward functions through individual
269 training with 5 different seeds and measured the quality of the reward functions by reporting the
270 average maximum evaluation score achieved at each policy checkpoint.

271 **Evaluation metrics.** In the Bidex tasks, the success of a task is determined by a binary 0-1 signal,
272 thus the evaluation score directly based on the success rate. In the Isaac tasks, the metric depends on
273 its goal. For example, the goal of the Ant task is to move as quickly as possible, so the metric is set
274 as the forward speed. See appendix for detailed task descriptions and evaluation scores.

275 **Method Implementation.** Since Eureka, Revolve, and RF-Agent optimize reward functions through
276 sampling, we ensure a fair comparison by setting the same total sampling number of reward functions.
277 In IsaacGym, we set the total limit to 80 and use the GPT-4o-mini-0718 and GPT-4o-0806 models
278 for LLM implementation[20]. For the Bidex tasks, to evaluate the search performance of the method
279 itself under complex control tasks, we increase the upper limit to 512 and use only the GPT-4o-mini
280 model. For Eureka and Revolve deployment details, please refer to the Appendix D. As for the
281 configuration of our RF-Agent, please refer to the Appendix C.

282 **5.4 Experiment Results**

283 **Advantages of RF-Agent in multi-class control tasks.** Tab.1 reports the absolute scores with
 284 standard deviation and average normalized performance of different methods on each IsaacGym
 285 task. We observe the following advantages: (1) Our method outperforms other LLM-based reward
 286 function design methods across various tasks and control types, demonstrating RF-Agent’s ability to
 287 leverage LLMs for reward function generation, with this advantage unaffected by changes in the LLM
 288 backbone model. (2) Even with a lightweight LLM backbone model 4o-mini, RF-Agent-generated
 289 reward functions maintain superiority over human expert methods in nearly all tasks, while Eureka
 290 and Revolve do not, particularly in locomotion tasks. (3) With a more powerful LLM backbone 4o,
 291 almost all methods achieve better reward functions within a limited number of training iterations.
 292 However, it is worth noting that Eureka experiences oscillations on FrankaCabinet, indicating its
 293 inability to generate effective reward functions with a lightweight model. Revolve fails to optimize
 294 performance on ShadowHand, while RF-Agent’s performance improvement remains stable across all
 295 tasks. These advantages highlight RF-Agent’s effective use of contextual reasoning and historical
 296 information, enabling superior search efficiency for high-performance reward functions.

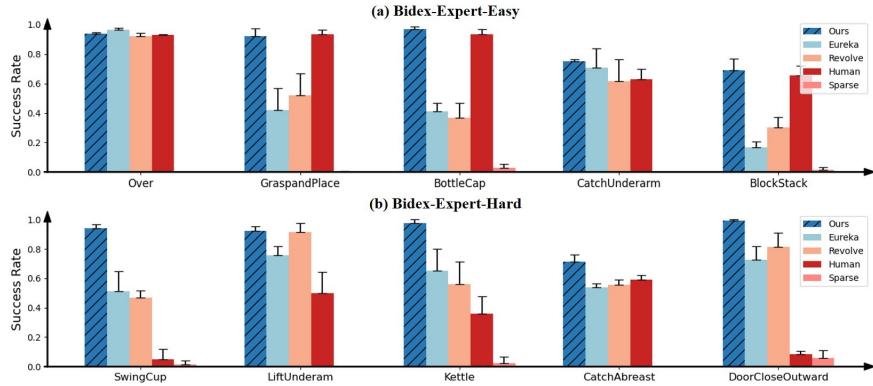


Figure 3: Success rates comparison with standard deviation upper bar on BiDex Expert-Easy/Hard.

297 **High Performance of RF-Agent under complex manipulation.** Fig.3 presents the results of
 298 different methods on BidexHands, with tasks divided into Expert-Easy and Expert-Hard groups
 299 based on human reward function outcomes for easier comparison. As shown, RF-Agent maintains a
 300 clear performance advantage over both human experts and other LLM-based reward function design
 301 methods, particularly on more complex tasks. Specifically, in the Expert-Easy tasks, where human
 302 reward functions have high success rates and smaller objects (*e.g.*, cubes, bottle caps) are manipulated,
 303 Eureka and Revolve fail to reach half of human performance on tasks like GraspAndPlace, BottleCap,
 304 and BlockStack. In contrast, RF-Agent matches or slightly exceeds human performance. In the
 305 Expert-Hard tasks, where human success rates are lower and larger objects (*e.g.*, kettles, doors) are
 306 manipulated, LLM-based methods perform all well, with RF-Agent showing an obvious advantage in
 307 most tasks. These results demonstrate that RF-Agent can design effective reward functions even in
 308 more complex task environments.

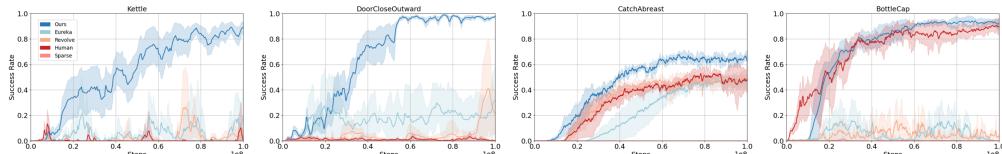


Figure 4: Success rates with exploration steps under reward functions by different methods.

309 **RF-Agent designs training-efficient reward functions.** Fig.4 shows the training curves for different
 310 reward functions on BiDexHands. Compared to Eureka and Revolve, Reward functions generated by
 311 RF-Agent can efficiently train the policy to converge to a higher success rate range, illustrating the
 312 high-quality reward function generation ability of RF-Agent. See Appendix F.1 for curves of others.

313 **The powerful search improvement capabilities of RF-Agent.** Fig.5 shows the average maximum
 314 scores achieved across tasks as reward functions are iteratively generated, reflecting the actual
 315 improvement of reward functions. In both Expert-Easy/Hard tasks, RF-Agent demonstrates high

316 optimization efficiency, significantly outperforming other methods. These results indicate that RF-
 317 Agent effectively addresses the issue of limited reward function improvement in complex control.
 318 We also provide the maximum reward achieved on each single task in Appendix F.2.

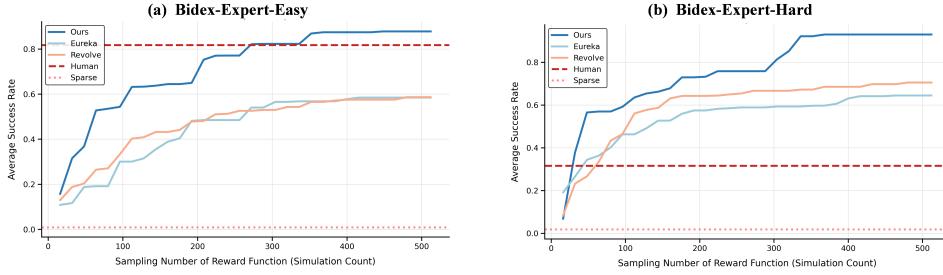


Figure 5: Reward function optimization performance with sampling counts.

319 5.5 Ablation Studies.

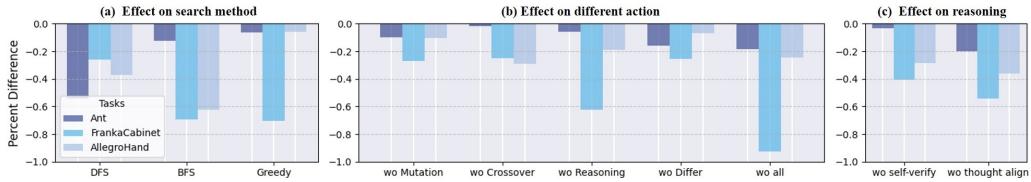


Figure 6: Ablations on the search method, action types, and reasoning component. We report the percent difference in performance compared to the complete RF-Agent. The experiment is built on Ant, FrankaCabinet, and AllegroHand to cover three task types with 4o-mini as the designer.

320 **Search Method.** To evaluate the impact of balancing exploration and exploitation on RF-Agent, we
 321 replace the improved UCT-based selection process with common search methods such as DFS, BFS,
 322 and Greedy. As shown in Fig 6.a, each method experienced significant performance loss on at least
 323 one task, highlighting the importance of balance for search efficiency.

324 **Action types.** To evaluate the impact of the actions designed in the expansion stage, we relatively
 325 replaced the different action types with a basic action, providing only the parent node state to generate
 326 new reward functions. In Fig 6.b, each action type holds its value, and removing all actions leads to a
 327 significant degradation. This demonstrates that the action design of RF-Agent effectively leverages
 328 historical information and LLM in-context learning ability to generate high-quality reward functions.

329 **Reasoning.** To evaluate the impact of LLM reasoning paradigm on RF-Agent, we remove self-verify
 330 and thought-align in Fig 6.c. The results show that the reasoning paradigm helps mitigate LLM
 331 hallucinations and provides effective evaluations, benefiting RF-Agent in generation and search.

332 6 Conclusion

333 This paper proposes RF-Agent, which automates the design of high-performance reward functions for
 334 low-level control tasks by combining language agents and tree search. By treating the reward function
 335 design process as a decision problem, RF-Agent integrates MCTS into the reward function design
 336 process, utilizing the multi-stage contextual reasoning of LLMs to improve search efficiency and
 337 leveraging diverse action design with historical information to make effective improvements to the
 338 reward function. The results demonstrate that RF-Agent can effectively generate high-performance
 339 reward functions for various and complex low-level control tasks, validating its effectiveness.

340 **Limitation.** Our RF-Agent, similar to Eureka and Revolve, belongs to the category of reward function
 341 design methods that utilize feedback and LLM. These methods are rather computationally expensive
 342 and time-consuming due to requiring multiple interactions with the LLM and repeated policy training.
 343 While RF-Agent achieves better generation results with relatively smaller LLM models, it does not
 344 reduce the need for multiple RL training iterations. Our future work will focus on reducing the
 345 number of RL training cycles while maintaining the effective iterative improvement.

346 **References**

- 347 [1] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning.
348 In *Proceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- 349 [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David,
350 Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not
351 as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- 352 [3] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob Mc-
353 Grew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning
354 dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20,
355 2020.
- 356 [4] Serena Booth, W Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi.
357 The perils of trial-and-error reward design: misdesign through overfitting and invalid task
358 specifications. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
359 pages 5920–5929, 2023.
- 360 [5] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling,
361 Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton.
362 A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence
363 and AI in games*, 4(1):1–43, 2012.
- 364 [6] Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-
365 Yves Oudeyer. Grounding large language models in interactive environments with online
366 reinforcement learning. In *International Conference on Machine Learning*, pages 3676–3713.
367 PMLR, 2023.
- 368 [7] Guillaume MJ B Chaslot, Mark HM Winands, and H Jaap van Den Herik. Parallel monte-carlo
369 tree search. In *Computers and Games: 6th International Conference, CG 2008, Beijing, China,
370 September 29-October 1, 2008. Proceedings 6*, pages 60–71. Springer, 2008.
- 371 [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
372 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
373 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 374 [9] Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong Feng, Jiechuan Jiang, Zongqing Lu,
375 Stephen McAleer, Hao Dong, Song-Chun Zhu, and Yaodong Yang. Towards human-level
376 bimanual dexterous manipulation with reinforcement learning. *Advances in Neural Information
377 Processing Systems*, 35:5150–5163, 2022.
- 378 [10] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
379 reinforcement learning from human preferences. *Advances in neural information processing
380 systems*, 30, 2017.
- 381 [11] Kenneth De Jong. Evolutionary computation: a unified approach. In *Proceedings of the Genetic
382 and Evolutionary Computation Conference Companion*, pages 373–388, 2017.
- 383 [12] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek
384 Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language
385 models. In *International Conference on Machine Learning*, pages 8657–8677. PMLR, 2023.
- 386 [13] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal
387 control via policy optimization. In *International conference on machine learning*, pages 49–58.
388 PMLR, 2016.
- 389 [14] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
390 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
391 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- 392 [15] Ankur Handa, Arthur Allshire, Viktor Makoviychuk, Aleksei Petrenko, Ritvik Singh, Jingzhou
 393 Liu, Denys Makoviichuk, Karl Van Wyk, Alexander Zhurkevich, Balakumar Sundaralingam,
 394 et al. Dextreme: Transfer of agile in-hand manipulation from simulation to reality. In *2023*
 395 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 5977–5984. IEEE,
 396 2023.
- 397 [16] Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhit-
 398 ing Hu. Reasoning with language model is planning with world model. *arXiv preprint*
 399 *arXiv:2305.14992*, 2023.
- 400 [17] RISHI HAZRA, Alkis Sygkounas, Andreas Persson, Amy Loutfi, and Pedro Zuidberg Dos
 401 Martires. REvolve: Reward evolution with large language models using human feedback. In *The Thirteenth International Conference on Learning Representations*, 2025.
- 403 [18] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
 404 *information processing systems*, 29, 2016.
- 405 [19] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qiang-
 406 long Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large
 407 language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on*
 408 *Information Systems*, 43(2):1–55, 2025.
- 409 [20] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
 410 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv*
 411 *preprint arXiv:2410.21276*, 2024.
- 412 [21] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward
 413 learning from human preferences and demonstrations in atari. *Advances in neural information*
 414 *processing systems*, 31, 2018.
- 415 [22] W Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward
 416 (mis) design for autonomous driving. *Artificial Intelligence*, 316:103829, 2023.
- 417 [23] Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with
 418 language models. *arXiv preprint arXiv:2303.00001*, 2023.
- 419 [24] Xinzhe Li. A survey on llm test-time compute via search: Tasks, llm profiling, search algorithms,
 420 and relevant frameworks. *arXiv preprint arXiv:2501.10069*, 2025.
- 421 [25] Yuxi Li. Deep reinforcement learning: An overview. *arXiv preprint arXiv:1701.07274*, 2017.
- 422 [26] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence,
 423 and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023*
 424 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE,
 425 2023.
- 426 [27] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh
 427 Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design
 428 via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- 429 [28] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 430 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
 431 with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594,
 432 2023.
- 433 [29] Denys Makoviichuk and Viktor Makoviychuk. rl-games: A high-performance framework for
 434 reinforcement learning. https://github.com/Denys88/rl_games, May 2021.
- 435 [30] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
 436 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High
 437 performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*,
 438 2021.

- 439 [31] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted
 440 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 442 [32] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad
 443 Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large
 444 language models. *arXiv preprint arXiv:2307.06435*, 2023.
- 445 [33] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transforma-
 446 tions: Theory and application to reward shaping. In *Icml*, volume 99, pages 278–287. Citeseer,
 447 1999.
- 448 [34] Jongjin Park, Younggyo Seo, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Surf:
 449 Semi-supervised reward learning with data augmentation for feedback-efficient preference-based
 450 reinforcement learning. *arXiv preprint arXiv:2203.10050*, 2022.
- 451 [35] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and
 452 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceed-
 453 ings of the 36th annual acm symposium on user interface software and technology*, pages 1–22,
 454 2023.
- 455 [36] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*.
 456 John Wiley & Sons, 2014.
- 457 [37] Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lyra Zhang, Fan Yang, and Mao Yang. Mutual
 458 reasoning makes smaller llms stronger problem-solvers. *arXiv preprint arXiv:2408.06195*,
 459 2024.
- 460 [38] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan,
 461 Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation
 462 models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- 463 [39] Nikita Rudin, David Hoeller, Marko Bjelonic, and Marco Hutter. Advanced skills by learning
 464 locomotion and local navigation end-to-end. In *2022 IEEE/RSJ International Conference on
 465 Intelligent Robots and Systems (IROS)*, pages 2497–2503. IEEE, 2022.
- 466 [40] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
 467 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 468 [41] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
 469 Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- 470 [42] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and
 471 Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive
 472 learning. *arXiv preprint arXiv:2010.03768*, 2020.
- 473 [43] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay,
 474 Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task
 475 plans using large language models. In *2023 IEEE International Conference on Robotics and
 476 Automation (ICRA)*, pages 11523–11530. IEEE, 2023.
- 477 [44] Satinder Singh, Richard L Lewis, and Andrew G Barto. Where do rewards come from. In
 478 *Proceedings of the annual conference of the cognitive science society*, pages 2601–2606.
 479 Cognitive Science Society, 2009.
- 480 [45] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1.
 481 MIT press Cambridge, 1998.
- 482 [46] Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True
 483 knowledge comes from practice: Aligning llms with embodied environments via reinforcement
 484 learning. *arXiv preprint arXiv:2401.14151*, 2024.
- 485 [47] Andrew Vince. A framework for the greedy algorithm. *Discrete Applied Mathematics*, 121(1-
 486 3):247–260, 2002.

- 487 [48] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
 488 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models.
 489 *arXiv preprint arXiv:2305.16291*, 2023.
- 490 [49] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen,
 491 Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous
 492 agents. *Frontiers of Computer Science*, 18(6):186345, 2024.
- 493 [50] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha
 494 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
 495 models. *arXiv preprint arXiv:2203.11171*, 2022.
- 496 [51] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani
 497 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large
 498 language models. *arXiv preprint arXiv:2206.07682*, 2022.
- 499 [52] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
 500 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
 501 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 502 [53] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse
 503 reinforcement learning. *arXiv preprint arXiv:1507.04888*, 2015.
- 504 [54] Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang,
 505 and Tao Yu. Text2reward: Automated dense reward function generation for reinforcement
 506 learning. In *International Conference on Learning Representations (ICLR)*, 2024 (07/05/2024–
 507 11/05/2024, Vienna, Austria), 2024.
- 508 [55] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
 509 Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint
 510 arXiv:2412.15115*, 2024.
- 511 [56] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
 512 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-
 513 vances in neural information processing systems*, 36:11809–11822, 2023.
- 514 [57] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan
 515 Cao. React: Synergizing reasoning and acting in language models. In *International Conference
 516 on Learning Representations (ICLR)*, 2023.
- 517 [58] Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez
 518 Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplík, et al. Language
 519 to rewards for robotic skill synthesis. *arXiv preprint arXiv:2306.08647*, 2023.
- 520 [59] Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu, Bingchao Wu, Bei Guan, Yongji Wang,
 521 and Jian-Guang Lou. Large language models meet nl2code: A survey. *arXiv preprint
 522 arXiv:2212.09420*, 2022.
- 523 [60] Yuwei Zeng, Yao Mu, and Lin Shao. Learning reward for robot skills using large language
 524 models via self-alignment. *arXiv preprint arXiv:2405.07162*, 2024.
- 525 [61] Chen Bo Calvin Zhang, Zhang-Wei Hong, Aldo Pacchiano, and Pulkit Agrawal. Orso: Ac-
 526 celerating reward design via online reward selection and policy optimization. *arXiv preprint
 527 arXiv:2410.13837*, 2024.
- 528 [62] Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge
 529 for large-scale task planning. In *Thirty-seventh Conference on Neural Information Processing
 530 Systems*, 2023.
- 531 [63] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang.
 532 Language agent tree search unifies reasoning acting and planning in language models. *arXiv
 533 preprint arXiv:2310.04406*, 2023.
- 534 [64] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy
 535 inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.

536 **A Environments**

537 In this section, we provide environment details. For each environment, we list its observation and
 538 action dimensions, the verbatim task description, and the task evaluation score function F .

539 **IsaacGym.** The general objective of the IsaacGym[30] environments is to simulate various robotic
 540 tasks in a controlled virtual setting. These tasks range from balancing and controlling movements to
 541 interacting with objects. Each environment has specific goals such as making an ant run as fast as
 542 possible, having a humanoid perform tasks with agility, or controlling a hand to spin an object to a
 543 target orientation. The environments involve challenges in movement, manipulation, and stability,
 544 with each task requiring precise control of the robot or agent, measured by a corresponding evaluation
 545 score function that encourages the agent to optimize its performance.

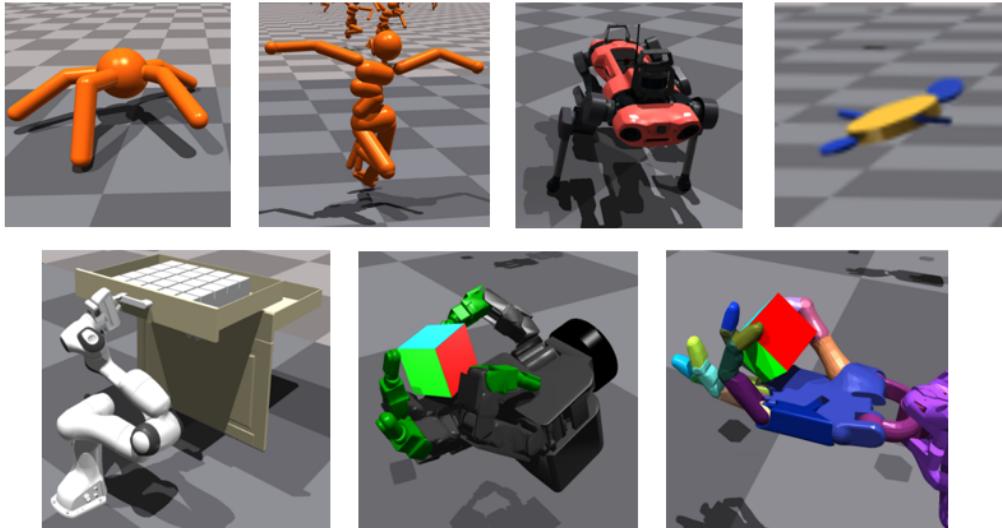


Figure 7: Examples of the IsaacGym environment.

IsaacGym Environments

Environment (obs dim, action dim)

Task description

Task evaluation score function F

Ant (60,8)

To make the ant run forward as fast as possible

`cur_dist - prev_dist`

Anymal (48, 12)

To make the quadruped follow randomly chosen x, y, and yaw target velocities

`- (linvel_error + angvel_error)`

Humanoid (108, 21)

To make the humanoid run as fast as possible

`cur_dist - prev_dist`

Quadcopter (21, 12)

To make the quadcopter reach and hover near a fixed position

`-cur_dist`

AllegroHand (88, 16)

To make the hand spin the object to a target orientation

```
number of consecutive successes where current success is 1[rot_dist < 0.1]
```

FrankaCabinet (23, 9)

To open the cabinet door

```
1[cabinet_pos > 0.39]
```

ShadowHand (211, 20)

To make the shadow hand spin the object to a target orientation

```
number of consecutive successes where current success is 1[rot_dist < 0.1]
```

546 **Bi-DexHands.** The Bi-DexHands[9] environment is designed to simulate dexterous manipulation
547 tasks involving two hands. The overall goal of the Bi-DexHand environment is to enable the agent
548 to perform various tasks that require coordinated manipulation using both hands. These tasks
549 are modeled to simulate real-world dexterous activities, such as passing objects between hands,
550 interacting with everyday tools, and performing precise actions requiring dual-hand cooperation.
551 Tasks in Bi-DexHands environment are designed to push the boundaries of agent control, testing their
552 ability to handle complex multi-step processes, balance, and precision, all of which are crucial in
553 dexterous manipulation.

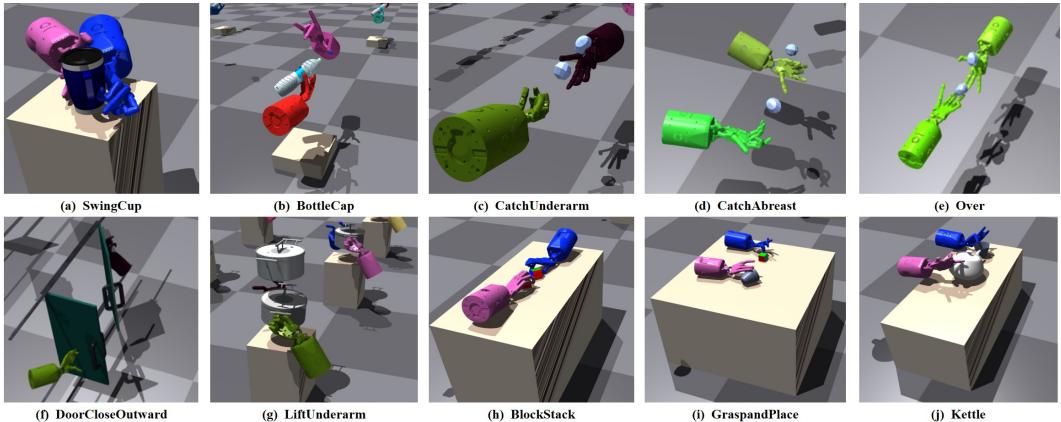


Figure 8: Examples of the Bi-DexHands environment.

Bi-DexHands Environments

Environment (obs dim, action dim)

Task description

Task evaluation score function F

Over (398, 40)

This class corresponds to the HandOver task. This environment consists of two shadow hands with palms facing up, opposite each other, and an object that needs to be passed. In the beginning, the object will fall randomly in the area of the shadow hand on the right side. Then the hand holds the object and passes the object to the other hand. Note that the base of the hand is fixed. More importantly, the hand which holds the object initially can not directly touch the target, nor can it directly roll the object to the other hand, so the object must be thrown up and stays in the air in the process

```
1[dist < 0.03]
```

GraspAndPlace (425, 52)

This class corresponds to the GraspAndPlace task. This environment consists of dual-hands, an object and a bucket that requires us to pick up the object and put it into the bucket

```
1[|block- bucket| < 0.2]
```

554 **BottleCap** (420, 52)

This class corresponds to the Bottle Cap task. This environment involves two hands and a bottle, we need to hold the bottle with one hand and open the bottle cap with the other hand. This skill requires the cooperation of two hands to ensure that the cap does not fall

555 1 [dist > 0.03]

556 **CatchUnderarm** (422, 52)

This class corresponds to the Catch Underarm task. In this task, two shadow hands with palms facing upwards are controlled to pass an object from one palm to the other. What makes it more difficult than the Hand over task is that the hands' translation and rotation degrees of freedom are no longer frozen but are added into the action space

557 1 [dist < 0.03]

558 **BlockStack** (428, 52)

This class corresponds to the Block Stack task. This environment involves dual hands and two blocks, and we need to stack the block as a tower

559 1 [goal_dist_1 < 0.07 and goal_dist_2 < 0.07 and 50 * (0.05- z_dist_1) > 1]

560 **SwingCup** (417, 52)

This class corresponds to the SwingCup task. This environment involves two hands and a dual handle cup, we need to use two hands to hold and swing the cup together

561 1 [rot_dist < 0.785]

562 **LiftUnderarm** (417, 52)

This class corresponds to the LiftUnderarm task. This environment requires grasping the pot handle with two hands and lifting the pot to the designated position. This environment is designed to simulate the scene of lift in daily life and is a practical skill

563 1 [dist < 0.05]

564 **Kettle** (417, 52)

This class corresponds to the PourWater task. This environment involves two hands, a kettle, and a bucket, we need to hold the kettle with one hand and the bucket with the other hand, and pour the water from the kettle into the bucket. In the practice task in Isaac Gym, we use many small balls to simulate the water

565 1 [|bucket- kettle_spout| < 0.05]

566 **CatchAbreast** (422, 52)

This class corresponds to the Catch Abreast task. This environment consists of two shadow hands placed side by side in the same direction and an object that needs to be passed. Compared with the previous environment which is more like passing objects between the hands of two people, this environment is designed to simulate the two hands of the same person passing objects, so different catch techniques are also required and require more hand translation and rotation techniques

567 1 [dist] < 0.03

568 **DoorCloseOutward** (417, 52)

This class corresponds to the DoorCloseOutward task. This environment also require a closed door to be opened and the door can only be pushed inward or initially open outward, but because they can't complete the task by simply pushing, which need to catch the handle by hand and then open or close it, so it is relatively difficult

569 1 [door_handle_dist < 0.5]

570 **B RF-Agent Details**

571 **B.1 Basic Prompts.**

572 Here, we present the fundamental prompts used across all LLM-based reward function design
573 methods, including the system prompt, environment information prompt, code format tip prompt.

574 The **system prompt** is as follows:

System Prompt

You are a reward engineer trying to write reward functions to solve reinforcement learning tasks as effective as possible.

Your goal is to write a reward function for the environment that will help the agent learn the task described in text.

Your reward function should use useful variables from the environment as inputs. As an example, the reward function signature can be:

```
@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor) -> Tuple[
    torch.Tensor, Dict[str, torch.Tensor]]:
    ...
    return reward, {}
```

Since the reward function will be decorated with `@torch.jit.script`, please make sure that the code is compatible with TorchScript (e.g., use torch tensor instead of numpy array).

Make sure any new tensor or variable you introduce is on the same device as the input tensors.

559

560 The **environment information prompt** with task **IsaacGym-Ant** is as follows:

Environment information Prompt with Task Ant

The task is: Ant. You need to make the ant run forward as fast as possible.

The Python environment is:

```
class Ant(VecTask):
    """Rest of the environment definition omitted."""
    def compute_observations(self):
        self.gym.refresh_dof_state_tensor(self.sim)
        self.gym.refresh_actor_root_state_tensor(self.sim)
        self.gym.refresh_force_sensor_tensor(self.sim)

        self.obs_buf[:, self.potentials[:, :], self.prev_potentials[:, :], self.up_vec
                    [:], self.heading_vec[:] = compute_ant_observations(
            self.obs_buf, self.root_states, self.targets, self.potentials,
            self.inv_start_rot, self.dof_pos, self.dof_vel,
            self.dof_limits_lower, self.dof_limits_upper, self.dof_vel_scale,
            self.vec_sensor_tensor, self.actions, self.dt, self.
            contact_force_scale,
            self.basis_vec0, self.basis_vec1, self.up_axis_idx)

    def compute_ant_observations(obs_buf, root_states, targets, potentials,
                                 inv_start_rot, dof_pos, dof_vel,
                                 dof_limits_lower, dof_limits_upper, dof_vel_scale,
                                 sensor_force_torques, actions, dt,
                                 contact_force_scale,
                                 basis_vec0, basis_vec1, up_axis_idx):
        # type: (Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, Tensor,
        #        Tensor, float, Tensor, float, float, Tensor, Tensor, int) ->
        # Tuple[Tensor, Tensor, Tensor, Tensor]
        torso_position = root_states[:, 0:3]
        torso_rotation = root_states[:, 3:7]
        velocity = root_states[:, 7:10]
        ang_velocity = root_states[:, 10:13]

        to_target = targets - torso_position
        to_target[:, 2] = 0.0

        prev_potentials_new = potentials.clone()
        potentials = -torch.norm(to_target, p=2, dim=-1) / dt

        torso_quat, up_proj, heading_proj, up_vec, heading_vec =
            compute_heading_and_up(
                torso_rotation, inv_start_rot, to_target, basis_vec0, basis_vec1, 2)
```

561

```

    vel_loc, angvel_loc, roll, pitch, yaw, angle_to_target = compute_rot(
        torso_quat, velocity, ang_velocity, targets, torso_position)

    dof_pos_scaled = unscale(dof_pos, dof_limits_lower, dof_limits_upper)

    obs = torch.cat((torso_position[:, up_axis_idx].view(-1, 1), vel_loc,
                     angvel_loc,
                     yaw.unsqueeze(-1), roll.unsqueeze(-1), angle_to_target,
                     unsqueeze(-1),
                     up_proj.unsqueeze(-1), heading_proj.unsqueeze(-1),
                     dof_pos_scaled,
                     dof_vel * dof_vel_scale, sensor_force_torques.view(-1, 24) *
                     contact_force_scale,
                     actions), dim=-1)

    return obs, potentials, prev_potentials_new, up_vec, heading_vec

```

562

563 The environment information prompt with task **BidexHands-Over** is as follows:

Environment information Prompt with Task Bidex-Over

The task is: This class corresponds to the HandOver task. This environment consists of two shadow hands with palms facing up, opposite each other, and an object that needs to be passed. In the beginning, the object will fall randomly in the area of the shadow hand on the right side. Then the hand holds the object and passes the object to the other hand. Note that the base of the hand is fixed. More importantly, the hand which holds the object initially can not directly touch the target, nor can it directly roll the object to the other hand, so the object must be thrown up and stays in the air in the process.

The Python environment is:

```

class ShadowHandOver(VecTask):
    """Rest of the environment definition omitted."""
    def compute_observations(self):
        self.gym.refresh_dof_state_tensor(self.sim)
        self.gym.refresh_actor_root_state_tensor(self.sim)
        self.gym.refresh_rigid_body_state_tensor(self.sim)
        self.gym.refresh_force_sensor_tensor(self.sim)
        self.gym.refresh_dof_force_tensor(self.sim)

        if self.obs_type in ["point_cloud"]:
            self.gym.render_all_camera_sensors(self.sim)
            self.gym.start_access_image_tensors(self.sim)

        self.object_pose = self.root_state_tensor[self.object_indices, 0:7]
        self.object_pos = self.root_state_tensor[self.object_indices, 0:3]
        self.object_rot = self.root_state_tensor[self.object_indices, 3:7]
        self.object_linvel = self.root_state_tensor[self.object_indices, 7:10]
        self.object_angvel = self.root_state_tensor[self.object_indices, 10:13]

        self.goal_pose = self.goal_states[:, 0:7]
        self.goal_pos = self.goal_states[:, 0:3]
        self.goal_rot = self.goal_states[:, 3:7]

        self.fingertip_state = self.rigid_body_states[:, self.fingertip_handles
            ][:, :, 0:13]
        self.fingertip_pos = self.rigid_body_states[:, self.fingertip_handles][:,
            :, 0:3]
        self.fingertip_another_state = self.rigid_body_states[:, self.
            fingertip_another_handles][:, :, 0:13]
        self.fingertip_another_pos = self.rigid_body_states[:, self.
            fingertip_another_handles][:, :, 0:3]

        if self.obs_type == "full_state":
            self.compute_full_state()
        elif self.obs_type == "point_cloud":
            self.compute_point_cloud_observation()

        if self.asymmetric_obs:
            self.compute_full_state(True)

```

564

565 The **code format tip prompt** is as follows:

Code format tip Prompt

Some helpful tips for writing the reward function code:

(1) You may find it helpful to normalize to a fixed reward range like -1,1 or 0,1 by applying transformations like `torch.exp()` to the overall reward or its reward components in order to have reasonable values for the optimization procedure. This does not mean that you should bound the rewards as wrong bounds might hinder the learning algorithm from achieving the optimal policy.

(2) If you choose to transform a reward component, then you must also introduce a temperature parameter inside the transformation function; this parameter must be a named variable in the reward function and it must not be an input variable. Each transformed reward component should have its own temperature variable.

(3) If the task has a goal state, it might be helpful to add a bonus to the reward function if the agent achieves the goal (up to some tolerance).

(4) Make sure the type of each input variable is correctly specified; a float input variable should not be specified as `torch.Tensor`.

(5) Make sure you don't give contradictory reward components.

(6) Most importantly, the reward code's input variables must contain only attributes of the provided environment class definition (namely, variables that have prefix `self.`). Under no circumstance can you introduce new input variables.

566

567 B.2 Feedback Example and Analysis Prompts.

568 Here we use the Ant task as an example to show the specific content of feedback. This feedback is
569 inherited from Eureka[27] and tracks the changes in a set of key variables during the training process,
570 including the key components of the designed reward function, task scores and episode length.

Feedback example on task Ant

`reward_forward_velocity:`

' -0.02', '1.07', '1.47', '1.90', '2.29', '2.62', '3.00', '3.48', '3.54', '3.67'

, Max: 3.73, Mean: 2.48, Min: -0.02
`reward_to_target:`

'0.00', '0.00', '0.00', '0.00', '0.00', '0.00', '0.00', '0.00', '0.00', '0.00'

, Max: 0.00, Mean: 0.00, Min: 0.00
`task_score:`

' -0.02', '1.08', '1.48', '1.90', '2.29', '2.60', '2.97', '3.45', '3.49', '3.62'

, Max: 3.67, Mean: 2.46, Min: -0.02
`episode_lengths:`

'59.19', '180.39', '285.69', '425.37', '519.28', '578.28', '634.42', '644.16', '633.70', '649.94'

, Max: 717.00, Mean: 494.84, Min: 59.19

571

572 Based on a given feedback, we provide some general tips for analyzing feedback.

Trained results analysis tip

(1) 'task_score' reflects the agent's actual task score or success rate after training under the current reward function.

573

- (2) If the values for a certain reward component are near identical throughout, then this means RL is not able to optimize this component as it is written.
 (3) If some reward components' magnitude is significantly larger or smaller, its value may not be conducive to policy learning.

574

575 The above feedback mechanism will be used in the implementation of Revolve to achieve automatic
 576 feedback of revolve in the experimental environment.

577 **B.3 Prompts for RF-Agent Actions.**

578 In RF-Agent, five different action types are set in the expansion stage, as well as the initial initialization
 579 action, as shown below:

580 The **initialization prompt** requires LLM to form a design idea based on the given environmental
 581 information and task description, and generate a reward function based on this idea.

Initialization prompt

First, **describe the design idea** and main steps of your reward function in one sentence. The description must be inside a brace outside the code implementation.

Next, **write a reward function** based on this idea. The output of the reward function should consist of two items:

- (1) the total reward,
- (2) a dictionary of each individual reward component.

The code output should be formatted as a python code string: "```python ... ```".
 Do not give additional explanations.

582

583 The **mutation a_{m1} prompt** requires the LLM to locally optimize its reward function based on the
 584 feedback information of the parent node. In particular, a_{m1} guides the LLM to modify the reward
 585 function composition mechanism.

Mutation a_{m1} prompt

I have one reward function with its design idea and code as follows.

Design Idea: {design_idea}
 Code: {reward_function}

We trained a RL policy using the provided reward function code and tracked the values of the individual components in the reward function as well as global policy metrics such as success rates and episode lengths after every {epoch_freq} epochs and the maximum, mean, minimum values encountered:

{trained_results}
 Analysis tips for trained results:
 {trained_result_analysis_tip}

Please create a new reward function that **has a different form but can be a modified version of the provided reward function**. The new reward function should have a higher task score. Try to introduce more novel idea and add or remove some basic reward components from the environment.

586

587 The **mutation a_{m2} prompt** requires the LLM to locally optimize its reward function based on the
 588 feedback information of the parent node. In particular, a_{m2} guides the LLM to adjust the parameter
 589 weights within the reward function.

Mutation a_{m2} prompt

I have one reward function with its design idea and code as follows.

Design Idea: {design_idea}

Code: {reward_function}

We trained a RL policy using the provided reward function code and tracked the values of the individual components in the reward function as well as global policy metrics such as success rates and episode lengths after every {epoch_freq} epochs and the maximum, mean, minimum values encountered:

```
{trained_results}  
Analysis tips for trained results:  
{trained_result_analysis_tip}
```

Please identify the reward function parameters and create a new reward function that has a different parameter settings compared to the provided version. The new reward function should have a higher task score.

590

591 The **crossover** a_{c3} **prompt** uses parent node information and nodes information in the elite set
592 to generate a new reward function by combining the reward component. After sorting the scores
593 corresponding to nodes in the set, the sampling of elite sets is randomly selected using the reciprocal
594 of sorting as the weight.

Crossover a_{c3} prompt

I have {nums} existing reward functions with their design ideas and codes as follows.

```
{reward_func_group}
```

We trained a RL policy using the provided reward function code and tracked the values of the individual components in the reward function as well as global policy metrics such as success rates and episode lengths after every {epoch_freq} epochs and the maximum, mean, minimum values encountered:

```
{trained_results}  
Analysis tips for trained results:  
{trained_result_analysis_tip}
```

Please create a new reward function inspired by those reward functions. Try to list the common ideas in those high score reward functions and combine high-performing reward components from different reward functions based on the data tracked. The new reward function should have a higher task score.

595

596 The **path reasoning** a_{r4} **prompt** uses parent node information and its ancestor information to
597 generate a new reward function by reasoning this optimization path.

Path Reasoning a_{r4} prompt

I have {nums} existing reward functions related to optimization in sequence with their design ideas and codes as follows.

```
{reward_func_group}
```

We trained a RL policy using the provided reward function code and tracked the values of the individual components in the reward function as well as global policy metrics such as

598

success rates and episode lengths after every {epoch_freq} epochs and the maximum, mean, minimum values encountered:

```
{trained_results}  
Analysis tips for trained results:  
{trained_result_analysis_tip}
```

Please create a new reward function inspired by all the above reward functions. **Try to reason the optimization path and list some ideas in those reward functions that are clearly helpful to a better improvement.** The new reward function should have a higher task score than any of them.

599

600 The **different thought a_{d5} prompt** uses parent node information and random nodes information
601 from the tree to generate a new reward function by requiring to design a new reward function with
602 different structures from these nodes.

Different Thought a_{d5} prompt

I have {nums} existing reward functions with their design ideas and codes as follows.

```
{reward_func_group}
```

We trained a RL policy using the provided reward function code and tracked the values of the individual components in the reward function as well as global policy metrics such as success rates and episode lengths after every {epoch_freq} epochs and the maximum, mean, minimum values encountered:

```
{trained_results}  
Analysis tips for trained results:  
{trained_result_analysis_tip}
```

Please create a new reward function that **has a totally different form from the given algorithms.** **Try generating codes with different structures, flows or algorithms.** Remember the new reward function should have a higher task score.

603

604 **B.4 Prompts for Thought-align and Self-verify.**

605 **Thought-align.** After the reward function is successfully executed, the RF-Agent resumaries
606 the reward function to eliminate the inconsistency between the original design idea and the reward
607 function expression caused by hallucinations or modifications.

Thought-align prompt

Following is the Design Idea of a reward function for the problem and the code for implementing the reward function.

Design Idea: {design_idea}
Code: {reward_function}

The content of the Design Idea cannot fully represent what the reward function has done informative. So, now you should re-describe the reward function using less than 4 sentences. Hint: You should reference the given Design Idea and highlight the most critical design ideas of the code. You can analyze the code to describe which reward components it contains, which observations they are calculated based on, what are the parameters and structure of the reward coupling process, and what special ideas are used.

608

609 **Self-verify.** In order to identify potential functions faster in the early selection stage, RF-Agent
610 appropriately added a self-verify score from the LLM as part of the selection criteria, and the impact
611 of this score will gradually weaken in the later stage.

Self-verify prompt

Be sure to remember the definition of the task. Additionally, I have an existing reward function with its design idea and code as follows.

Design Idea: {design_idea}
Code: {reward_function}

Now, based on the task definition, imagine how an expert-level strategy completes the task, describe the execution planning and motion process of the expert strategy for the task, and then evaluate the given reward functions according to your imagined description process, and judge the similarity of the given reward function with the imagined process, with the numerical value limited to [-1,1].

Finally, return the value enclosed in square brackets, such as [0.5].

612

613 C Training Details

614 **Hyper-parameters.** As shown in the paper, we deploy our RF-Agent with leaf node parallel
615 scheme[7] in the expansion stage, we set this parallel number to 8. Thus, we configure the actions as
616 [2, 2, 2, 1, 1] per expansion to make the ratio of nodes that utilize local and global information at 1 : 1.
617 The initial value of λ is set to 0.4 across all tasks, v_{self} is constrained within the range of $[-1, 1]$ in
618 order to have a more obvious distinction after softmax, k is randomly sampled from [2, 4] to control
619 the number of nodes sampled in a_{c3} , a_{r4} and a_{d5} , and η is set to 0.7.

620 **Deployment resources.** We deployed RF-Agent on a 4 Nvidia Geforce RTX3090 cards with 128
621 core CPUs and 256GiB memory server. On issacgym, the average time consumed by each different
622 task was 9 hours; on Bi-DexHands, the average time consumed by each task was 40 hours.

623 **Other evaluation details.** During the evaluation, we tested on 5 seeds different from those during
624 the search. The number of training steps in the environment was consistent with the default training
625 steps in the IssacGym and Bi-DexHands environments. For example, the uniform number of steps on
626 Bi-DexHands is 1e8.

627 D Baseline Details

628 **Eureka.** Eureka[27] has been tested on IsaacGym and Bi-Dexhands using the GPT-4 model as the
629 reward function generator. Given the high cost of GPT-4 and its reduced relevance as a mainstream
630 choice, we used 4o and 4o-mini models for our experiments. Therefore, we also re-conducted
631 experiments with Eureka using these two models. The parameter configuration for Eureka remains
632 consistent with its default settings, maintaining 16 samples per iteration.

633 **Revolve.** Since Revolve[17] has not been previously tested on IsaacGym and Bi-Dexhands, we
634 re-implemented the experiments. We retained its algorithm and evolutionary operation prompts, with
635 the main changes being in the system prompt for the environment and replacing feedback information
636 with automated environment feedback. These two types of prompts are consistent with those used in
637 Eureka and RF-Agent. The parameter configuration for Revolve remains consistent with its default
638 settings, maintaining 16 samples per iteration during population evolution.

639 E Cost Comparison

640 **Training resource consumption.** Eureka, Revolve and Ours RF-Agent are using the same total
641 number of reward functions to training the policy, thus the training cost is almost the same.

642 **Storage Usage.** RF-Agent needs to store historical information of the entire tree, Revolve needs to
 643 store historical information of the population, and Eureka hardly needs to store historical information,
 644 but these additional historical information is composed of strings, whose order of magnitude is at the
 645 MB level.

646 **Token consumption.** The main consumption of each method is concentrated on the process of
 647 generating reward functions in combination with environmental information and instructions. RF-
 648 Agent also has thought-align and self-verify processes compared to Eureka and Revolve, but the
 649 token used in these two processes is much lower than the reward function generation process.

650 Overall, RF-Agent uses slightly more storage space and token consumption, but experimental results
 651 show that this slight sacrifice is worth it.

652 F Additional Results

653 F.1 Training Curves in Bi-DexHands.

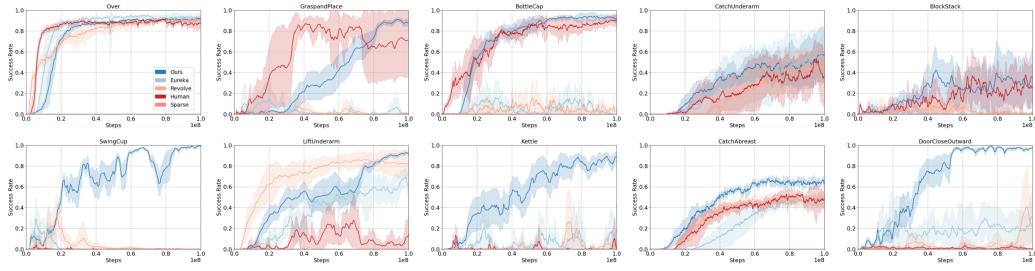


Figure 9: Success rates with exploration steps under reward functions by different methods on BiDexHands.

654 In Fig.9, we show all the training curve results of the Expert-Easy and Expert-Hard groups on the
 655 selected Bi-DexHands. RF-Agent can maintain its dominant performance in most tasks.

656 F.2 Reward Function Optimization Performance.

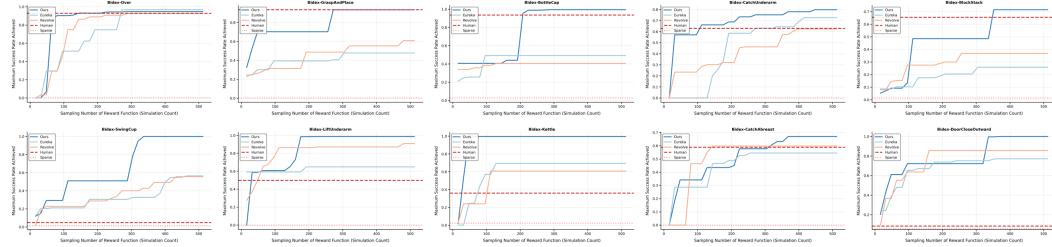


Figure 10: Reward function optimization performance with sampling counts in Bi-DexHands.

657 In Fig.10, we show the maximum score obtained on each task as the number of reward function
 658 samples increases. The results show that our RF-Agent has better iterative optimization capabilities
 659 on most tasks.

660 F.3 More Optimization Path Examples of RF-Agent

661 We show the optimization paths of the nodes that end up with the best performance on two Bi-
 662 DexHands tasks, as shown in Fig.11. We found that the optimization paths on both tasks contain
 663 certain global actions, which shows that it is important to optimize reward functions using global
 664 historical information. In addition, the optimization paths on CatchAbreast further illustrate that
 665 those weak reward functions should not be easily abandoned, and mixing them with reward functions
 666 on other paths can also obtain performance advantage nodes.

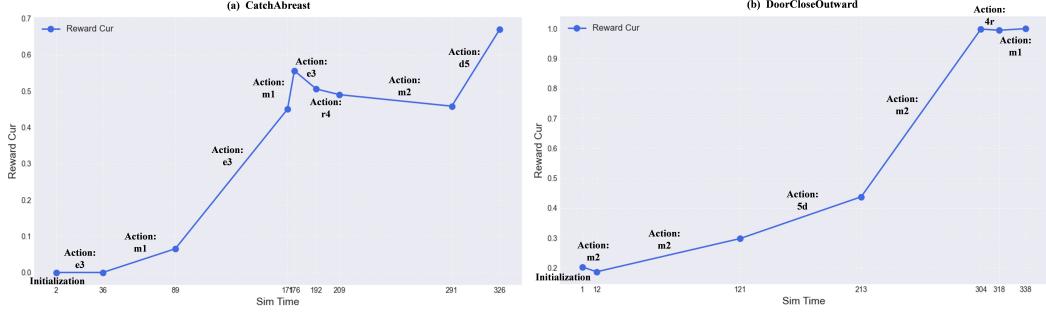


Figure 11: Examples of the best node growth path.

667 G Algorithm of RF-Agent

668 Algorithm 1 provides a pseudo-code for the proposed RF-Agent method, which can be combined
 669 with Fig.2 to further familiarize yourself with the entire RF-Agent process.

Algorithm 1: Reward Function Design via Language Agent Tree Search (RF-Agent)

Input: score Function F , reward designer G , the number of initial nodes N_I , max evaluation times N , action set $[a_{m1}, a_{m2}, a_{c3}, a_{r4}, a_{d5}]$, counts of action $[n_{m1}, n_{m2}, n_{c3}, n_{d4}, n_{d5}]$, UCT initial balance parameter λ_0

1 // **Initialization Stage**
 2 Initialize a virtual root node n_{root}
 3 Set $t \leftarrow 0$
 4 Initialization N_I nodes with designer G and link all the N_I nodes to the n_{root} ; // Eq.1
 5 Evaluating F after training policy π with R from N_I nodes
 6 // **MCTS Search**
 7 **for** $t \leq N$ **do**
 8 $\lambda \leftarrow \lambda_0 \cdot \frac{N-t}{N}$
 9 // **MCTS Selection Stage**
 10 $s \leftarrow n_{root}$
 11 **while** s is not a leaf node **do**
 12 $s \leftarrow \text{argmax}_{s' \in \text{Child}(s)} \text{UCT}(s')$; // Eq.2 with λ
 13 **end while**
 14 // **MCTS Expansion Stage**
 15 **for** a in action set **do**
 16 | expand node s with action a and designer G in n_a times ; // Eq.1 with l_{action}
 17 **end**
 18 // **MCTS Simulation Stage**
 19 Evaluating F after training policy π with R from newly generated nodes
 20 // **MCTS Backpropagation Stage**
 21 Evaluate the self-verify score v_{self} of newly generated nodes.
 22 **while** s is not the root node n_{root} **do**
 23 | Update $Q(s)$ and $N(s)$; // Eq.3
 24 **end while**
 25 $t \leftarrow t + \sum_{a \in \text{action set}} n_a$
 26 **end**
 27 **Return:** The node s^* with the best reward function R^*

670 H Licenses

671 The licenses and URL of baselines, benchmarks and deployment algorithm are listed in 4.

Resources	Type	License	URL
Eureka[27]	Codes for Baseline	MIT	https://github.com/eureka-research/Eureka
Revolve[17]	Codes for Baseline	Available Online	https://github.com/RishiHazra/Revolve/tree/main
IsaacGym[30]	Codes for Benchmark	BSD 3-Clause License	https://github.com/isaac-sim/IsaacGymEnvs
Bi-DexHands[9]	Codes for Benchmark	Apache 2.0	https://github.com/PKU-MARL/DexterousHands
RL-games[29]	Codes for Algorithm	MIT	https://github.com/Denys88/rl_games

Table 4: The licenses and URL of baselines , benchmarks and deployment algorithm.

672 I RF-Agent Reward Examples

673 I.1 Reward after Diverse Actions Examples

674 We randomly selected a task and a corresponding node to show what the tendency of the new reward
 675 function generated by RF-Agent under different actions. Here we select the node with RF-Agent
 676 developed on CatchAbreast to the path "root-0i-e3-m1-e3" as the current node. First, the reward
 677 function of the node is displayed, and then the reward function generated by different actions under
 678 the node is displayed in turn.

Current reward function on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
    left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.
    Tensor, Dict[str, torch.Tensor]]:
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=-1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=-1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=-1)

    # Reward for being close to the goal
    distance_reward = torch.clamp(1.0 - (distance_to_goal / 0.5), min=0, max=1)

    # Reward for successfully passing the object, bonus for being within a
    # proximity threshold
    success_bonus = torch.where(distance_to_goal < 0.1, torch.tensor(1.0, device=
        object_pos.device), torch.tensor(0.0, device=object_pos.device))

    # Hand balance reward penalizing uneven movement
    hand_balance_reward = -torch.abs(distance_to_left_hand -
        distance_to_right_hand)

    # Proximity reward to encourage hands to efficiently reach the object
    combined_hand_distance = distance_to_left_hand + distance_to_right_hand
    proximity_reward = torch.clamp(1.0 - (combined_hand_distance / 0.8), min=0,
        max=1)

    # Smoothness penalty to encourage fluid hand movements
    hand_distance_between_hands = torch.norm(left_hand_pos - right_hand_pos, dim
        =-1)
    smooth_movement_penalty = 0.01 * hand_distance_between_hands

    # Total reward calculation
    total_reward = distance_reward + success_bonus + hand_balance_reward +
        proximity_reward - smooth_movement_penalty

    # Reward components dictionary
    reward_components = {
        "distance_reward": distance_reward,
        "success_bonus": success_bonus,
        "hand_balance_reward": hand_balance_reward,
        "proximity_reward": proximity_reward,
        "smooth_movement_penalty": smooth_movement_penalty
    }

    return total_reward, reward_components

```

679

680 **Mutation** a_{m1} .

Reward function after a_{m1} on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
    left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.
    Tensor, Dict[str, torch.Tensor]]:
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=-1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=-1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=-1)

    # Reward for being very close to the goal with increased distance tolerance
    distance_reward = torch.clamp(1.0 - (distance_to_goal / 0.6), min=0, max=1)

    # Success bonus increased for being within a tighter proximity threshold
    success_bonus = torch.where(distance_to_goal < 0.05, torch.tensor(1.0, device
        =object_pos.device), torch.tensor(0.0, device=object_pos.device))

    # Balance reward penalizing uneven movement but normalizing based on hand
    # distance
    hand_balance_penalty = -torch.abs(distance_to_left_hand -
        distance_to_right_hand)
    normalized_hand_distance = torch.clamp((distance_to_left_hand +
        distance_to_right_hand) / (0.6 + 0.6), min=0, max=1)
    hand_balance_reward = hand_balance_penalty * normalized_hand_distance

    # Proximity reward encouraging efficient hand placement
    combined_hand_distance = distance_to_left_hand + distance_to_right_hand
    proximity_reward = torch.clamp(1.0 - (combined_hand_distance / 0.75), min=0,
        max=1)

    # Smoothness bonus for transitions encouraging fluid hand movements
    hand_distance_between_hands = torch.norm(left_hand_pos - right_hand_pos, dim
        =-1)
    smooth_movement_bonus = torch.clamp(1.0 - (hand_distance_between_hands / 0.1),
        min=0, max=1)

    # Total reward calculation
    total_reward = (distance_reward + success_bonus + hand_balance_reward +
        proximity_reward + smooth_movement_bonus) / 1.5

    # Reward components dictionary
    reward_components = {
        "distance_reward": distance_reward,
        "success_bonus": success_bonus,
        "hand_balance_reward": hand_balance_reward,
        "proximity_reward": proximity_reward,
        "smooth_movement_bonus": smooth_movement_bonus
    }

    return total_reward, reward_components

```

681

682 where the key change is adding the normalizing mechanism to "hand_balance_penalty"

683 **Mutation** a_{m2} .

Reward function after a_{m2} on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
    left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.
    Tensor, Dict[str, torch.Tensor]]:
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=-1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=-1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=-1)

    # Reward for being close to the goal with a larger range
    distance_reward = torch.clamp(1.0 - (distance_to_goal / 0.3), min=0, max=1)

    # Larger success bonus for passing the object
    success_bonus = torch.where(distance_to_goal < 0.05, torch.tensor(2.0, device
        =object_pos.device), torch.tensor(0.0, device=object_pos.device))

    # Hand balance reward maintaining a moderate penalty

```

684

```

    hand_balance_reward = -0.5 * torch.abs(distance_to_left_hand -
                                             distance_to_right_hand)

    # Enhanced proximity reward encouraging efficient hand placement
    combined_hand_distance = distance_to_left_hand + distance_to_right_hand
    proximity_reward = torch.clamp(1.0 - (combined_hand_distance / 0.6), min=0,
                                   max=1)

    # Smoothness penalty slightly adjusted to minimize harsh penalties
    hand_distance_between_hands = torch.norm(left_hand_pos - right_hand_pos, dim
                                              = 1)
    smooth_movement_penalty = 0.005 * hand_distance_between_hands

    # Total reward calculation
    total_reward = distance_reward + success_bonus + hand_balance_reward +
                  proximity_reward - smooth_movement_penalty

    # Reward components dictionary
    reward_components = {
        "distance_reward": distance_reward,
        "success_bonus": success_bonus,
        "hand_balance_reward": hand_balance_reward,
        "proximity_reward": proximity_reward,
        "smooth_movement_penalty": smooth_movement_penalty
    }

    return total_reward, reward_components

```

685

686 where the key change is change the parameters in some reward components.

687 **Crossover a_{c3} .**

Reward function after a_{c3} on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
                  left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.
Tensor, Dict[str, torch.Tensor]]:
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=1)

    goal_reward = torch.clamp(1.0 - (distance_to_goal / 0.5), min=0, max=1) # Proximity to goal
    success_bonus = torch.where(distance_to_goal < 0.1, torch.tensor(2.0, device=
        object_pos.device), torch.tensor(0.0, device=object_pos.device)) # Success bonus

    hand_balance_reward = -torch.abs(distance_to_left_hand -
                                     distance_to_right_hand) # Penalizes uneven movement between hands
    hand_coordination_reward = torch.exp(-torch.norm(left_hand_pos -
                                                     right_hand_pos, dim=1)) # Encourages closer hand positions
    proximity_reward = -0.5 * (distance_to_left_hand + distance_to_right_hand) # Maintain proximity to the object

    total_reward = goal_reward + success_bonus + hand_balance_reward +
                  hand_coordination_reward + proximity_reward

    reward_components = {
        "goal_reward": goal_reward,
        "success_bonus": success_bonus,
        "hand_balance_reward": hand_balance_reward,
        "hand_coordination_reward": hand_coordination_reward,
        "proximity_reward": proximity_reward
    }

    return total_reward, reward_components

```

688

689 where the key change is to delete "combined_hand_distance" and "smooth_movement_penalty" from
690 the initial reward function while keeping other high-performing components.

691 **Path Reasoning a_{r4} .**

Reward function after a_{r4} on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
                  left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
    # Constants
    distance_threshold = 0.1
    success_bonus = 1.0
    smooth_movement_penalty = 0.01

    # Distance metrics
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=1)

    # Reward components
    distance_reward = torch.clamp(1.0 - (distance_to_goal / 0.5), min=0, max=1)
    # Reward for being close to the goal

    # Reward for successfully passing the object
    successful_pass = (distance_to_goal < distance_threshold).type(object_pos.dtype)
    reward_successful_pass = successful_pass * success_bonus

    # Hand balance reward to encourage synchronized movements
    hand_balance_reward = -torch.abs(distance_to_left_hand -
                                      distance_to_right_hand)

    # Proximity reward for effectively reaching the object
    combined_hand_distance = distance_to_left_hand + distance_to_right_hand
    proximity_reward = torch.clamp(1.0 - (combined_hand_distance / 0.8), min=0,
                                   max=1)

    # Smoothness penalty for fluid hand movements
    hand_distance_between_hands = torch.norm(left_hand_pos - right_hand_pos, dim=1)
    smooth_movement_penalty_reward = -smooth_movement_penalty *
                                      hand_distance_between_hands

    # Total reward calculation
    total_reward = distance_reward + reward_successful_pass + hand_balance_reward
    + proximity_reward + smooth_movement_penalty_reward

    # Reward components dictionary
    reward_components = {
        "distance_reward": distance_reward,
        "reward_successful_pass": reward_successful_pass,
        "hand_balance_reward": hand_balance_reward,
        "proximity_reward": proximity_reward,
        "smooth_movement_penalty_reward": smooth_movement_penalty_reward
    }

    return total_reward, reward_components

```

692

693 where the key change is that the coefficients in the original reward function have been redefined and
694 some nested methods have been changed.

695 **Different Thought** a_{d5} .

Reward function after a_{r4} on CatchAbreast.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
                  left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor, object_linvel: torch.Tensor) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
    # Design Idea: Utilize a nested reward structure focusing on progressive
    # subtasks of grasping and passing the object, rewarding successful
    # transitions and smooth trajectories.

    # Define parameters
    goal_tolerance = 0.1
    grasp_radius = 0.05

```

696

```

smoothness_weight = 0.2

# Calculate distances for nested rewards
distance_to_goal = torch.norm(object_pos - goal_pos, dim=-1)
distance_to_left_hand = torch.norm(object_pos - left_hand_pos, dim=-1)
distance_to_right_hand = torch.norm(object_pos - right_hand_pos, dim=-1)

# Nested reward for successful grasping (within grasp radius)
grasp_reward = torch.where(distance_to_left_hand < grasp_radius, torch.tensor(0.5, device=object_pos.device), torch.zeros_like(distance_to_left_hand))
grasp_reward += torch.where(distance_to_right_hand < grasp_radius, torch.tensor(0.5, device=object_pos.device), torch.zeros_like(distance_to_right_hand))

# Goal reward
goal_reward = torch.where(distance_to_goal < goal_tolerance, torch.tensor(1.0, device=object_pos.device), torch.zeros_like(distance_to_goal))

# Smoothness reward calculated as the negative exponential of object linear velocity
smoothness_reward = -torch.exp(-smoothness_weight * torch.norm(object_linvel, dim=-1))

# Total reward combining all components
total_reward = smoothness_reward + grasp_reward + goal_reward

# Reward components
reward_components = {
    "grasp_reward": grasp_reward,
    "goal_reward": goal_reward,
    "smoothness_reward": smoothness_reward
}

return total_reward, reward_components

```

697

698 where the change is to think about new reward forms, introducing a grasping angle and nesting it into
699 the distance-related rewards.

700 **I.2 Design Thought after Alignment**

701 Similarly, the current node "root-0i-e3-m1-e3" in the previous section is an example. The design
702 thought during the generation process is:

Before thought-align

The reward function is designed to optimize agent performance in the Catch Abreast task by rewarding careful proximity of the hands to the object, effective synchronization of hand movements, and incentivizing reaching the goal with a success bonus, while penalizing excessive movement to ensure fluid hand coordination.

703

704 After thought alignment:

After thought-align

The reward function is designed to enhance agent performance in the Catch Abreast task by incorporating multiple components that prioritize hand coordination and goal achievement. It evaluates proximity to the goal and the object, providing rewards for closeness and a success bonus for passing the object within a defined threshold. Additionally, it incorporates a hand balance reward to penalize uneven movements between the hands, a proximity reward to encourage efficient hand placement, and a smoothness penalty to promote fluid motion. The total reward is calculated by summing these components while ensuring that fluid hand coordination is incentivized.

705

706 Obviously, the design ideas after thinking alignment can more comprehensively summarize the
707 corresponding reward function.

708 **I.3 RF-Agent and Baselines Reward Comparison**

709 Here we list some best reward function from different methods.

710 **IsaacGym-Ant.**

Our RF-Agent.

```

@torch.jit.script
def compute_reward(root_states: torch.Tensor, heading_vec: torch.Tensor, up_vec:
    torch.Tensor) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
    # Extract the forward velocity in the x-direction
    velocity = root_states[:, 7:10]
    forward_velocity = velocity[:, 0]

    # Weight for different reward components
    heading_weight = 1.5
    upright_weight = 0.5
    y_z_penalty_weight = 0.3

    # Cosine similarity reward for maintaining the desired heading direction (x-
        axis)
    desired_heading = torch.tensor([1.0, 0.0, 0.0], device=heading_vec.device).
        expand_as(heading_vec)
    heading_cos_sim = torch.nn.functional.cosine_similarity(desired_heading,
        heading_vec, dim=-1)
    heading_reward = heading_weight * (heading_cos_sim * forward_velocity)

    # Stability penalty for velocity in y and z directions
    y_z_velocity_penalty = y_z_penalty_weight * torch.norm(velocity[:, 1:], p=2,
        dim=-1)

    # Reward for being upright, using the z-component of the up_vec to measure
        uprightness
    upright_bonus_temp = 2.0
    upright_bonus = upright_weight * torch.exp(upright_bonus_temp * (up_vec[:, 2]
        - 1.0))

    # Total reward is a combination of forward movement, heading alignment, and
        uprightness
    total_reward = forward_velocity + heading_reward + upright_bonus -
        y_z_velocity_penalty

    # Construct reward components dictionary
    reward_components = {
        "forward_velocity": forward_velocity,
        "heading_reward": heading_reward,
        "upright_bonus": upright_bonus,
        "y_z_velocity_penalty": -y_z_velocity_penalty
    }

    return total_reward, reward_components

```

711

Eureka.

```

@torch.jit.script
def compute_reward(root_states: torch.Tensor, potentials: torch.Tensor,
    prev_potentials: torch.Tensor, actions: torch.Tensor,
    dof_vel: torch.Tensor, dof_limits_lower: torch.Tensor,
    dof_limits_upper: torch.Tensor) -> Tuple[torch.Tensor,
        Dict[str, torch.Tensor]]:
    # Extract relevant state information
    velocity = root_states[:, 7:10] # Linear velocity
    torso_position = root_states[:, 0:3] # Position of the torso

    # Reward component to encourage forward velocity
    forward_velocity_reward = velocity[:, 0] # Assuming forward is along the x-
        axis

    # Modify target closeness reward to reflect more sensitivity
    target_closeness_reward_temperate = 0.2
    target_closeness_reward = prev_potentials - potentials

```

712

```

exp_target_closeness_reward = torch.exp(target_closeness_reward_temperate *
target_closeness_reward)

# Re-examine dof velocity penalty for more sensitivity
dof_vel_penalty_temperate = 0.02
dof_vel_penalty = torch.sum(torch.abs(dof_vel), dim=-1)
exp_dof_vel_penalty = torch.exp(-dof_vel_penalty_temperate * dof_vel_penalty)

# Actions penalty adjust temperature for greater distinction
actions_penalty_temperate = 0.005
actions_penalty = torch.sum(torch.abs(actions), dim=-1)
exp_actions_penalty = torch.exp(-actions_penalty_temperate * actions_penalty)

# Total normalized reward
total_reward = forward_velocity_reward + 0.5 * exp_target_closeness_reward +
0.1 * exp_dof_vel_penalty + 0.1 * exp_actions_penalty

# Compile individual components into a dictionary
reward_components = {
    "forward_velocity_reward": forward_velocity_reward,
    "exp_target_closeness_reward": exp_target_closeness_reward,
    "exp_dof_vel_penalty": exp_dof_vel_penalty,
    "exp_actions_penalty": exp_actions_penalty,
}

return total_reward, reward_components

```

713

Revolve.

```

@torch.jit.script
def compute_reward(root_states: torch.Tensor, targets: torch.Tensor, potentials: torch.Tensor,
                  prev_potentials: torch.Tensor, up_vec: torch.Tensor) -> Tuple[
    torch.Tensor, Dict[str, torch.Tensor]]:
    # Extract necessary components
    velocity = root_states[:, 7:10]

    # Forward speed reward component
    forward_speed_reward = velocity[:, 0]  # Forward speed along x-axis

    # Forward potential difference reward component
    forward_potential_diff = potentials - prev_potentials

    # Upright reward component
    up_vec_goal = torch.tensor([0.0, 0.0, 1.0], device=up_vec.device).unsqueeze(
        0).expand_as(up_vec)
    upright_reward = torch.sum(up_vec * up_vec_goal, dim=-1)

    # Set temperatures for exponential scaling
    forward_speed_temp = 0.35
    forward_potential_temp = 0.4
    upright_temp = 1.0  # Keeping original scale

    # Exponentially scale the rewards
    scaled_forward_speed_reward = forward_speed_temp * torch.exp(
        forward_speed_reward * forward_speed_temp)
    scaled_forward_potential_diff = forward_potential_temp * torch.exp(
        forward_potential_diff * forward_potential_temp)
    scaled_upright_reward = torch.exp(upright_temp * upright_reward)

    # Total reward calculation
    total_reward = scaled_forward_speed_reward + scaled_forward_potential_diff +
    scaled_upright_reward

    # Reward components dictionary
    reward_components = {
        "forward_speed_reward": scaled_forward_speed_reward,
        "forward_potential_diff_reward": scaled_forward_potential_diff,
        "upright_reward": scaled_upright_reward
    }

    return total_reward, reward_components

```

714

715 **Bidex-CatchAbreast.**

Our RF-Agent.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
    left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor) -> Tuple[torch.
    Tensor, Dict[str, torch.Tensor]]:
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=1)
    distance_to_left_hand = torch.norm(left_hand_pos - object_pos, dim=1)
    distance_to_right_hand = torch.norm(right_hand_pos - object_pos, dim=1)

    # Reward for proximity to the goal with increased tolerance
    distance_reward = torch.clamp(1.0 - (distance_to_goal / 0.5), min=0, max=1)

    # Success bonus for accurate object passing with tighter threshold
    success_bonus = torch.where(distance_to_goal < 0.03, torch.tensor(1.0, device
        =object_pos.device), torch.tensor(0.0, device=object_pos.device))

    # Hand balance reward incentivizing even distribution of effort between hands
    hand_balance_reward = -0.5 * torch.abs(distance_to_left_hand -
        distance_to_right_hand) * torch.clamp(1.0 - (distance_to_left_hand +
        distance_to_right_hand) / 1.5, min=0, max=1)

    # Proximity reward to encourage both hands to be near the object
    combined_hand_distance = distance_to_left_hand + distance_to_right_hand
    proximity_reward = torch.clamp(1.0 - (combined_hand_distance / 0.75), min=0,
        max=1)

    # Smoothness bonus to promote fluid hand movements
    hand_distance_between_hands = torch.norm(left_hand_pos - right_hand_pos, dim
        =1)
    smooth_movement_bonus = torch.clamp(1.0 - (hand_distance_between_hands / 0.1),
        min=0, max=1)

    # Total reward calculation
    total_reward = (distance_reward + success_bonus + hand_balance_reward +
        proximity_reward + smooth_movement_bonus) / 2.0

    # Reward components dictionary
    reward_components = {
        "distance_reward": distance_reward,
        "success_bonus": success_bonus,
        "hand_balance_reward": hand_balance_reward,
        "proximity_reward": proximity_reward,
        "smooth_movement_bonus": smooth_movement_bonus
    }

    return total_reward, reward_components

```

716

Eureka.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
    object_linvel: torch.Tensor) -> Tuple[torch.Tensor, Dict[str,
    torch.Tensor]]:

    # Constants for temperature adjustments
    distance_temp = 5.0 # Increased temperature sensitivity for distance
    velocity_temp = 1.0 # Reduced scaling for velocity

    # Compute distance to goal
    distance_to_goal = torch.norm(object_pos - goal_pos, dim=1)
    # Negative distance to promote minimization
    distance_reward = torch.exp(-distance_temp * distance_to_goal) # Exponential
        decay for distance

    # Improved velocity reward with penalty for slow speeds
    speed_threshold = 0.1 # threshold for penalty
    velocity_magnitude = torch.norm(object_linvel, dim=1)
    # Penalty for low speeds to encourage proper movement
    velocity_penalty = torch.where(velocity_magnitude < speed_threshold, -1.0 * (
        speed_threshold - velocity_magnitude), torch.zeros_like(
        velocity_magnitude))

```

717

```

# Positive reward for aligned velocity towards the goal direction
direction_vector = goal_pos - object_pos # Direction vector from object to
                                         goal
direction_norm = torch.norm(direction_vector, dim=1, keepdim=True) + 1e-6
desired_velocity = direction_vector / direction_norm # Normalized direction
aligned_velocity = torch.sum(object_linvel * desired_velocity, dim=1) # Dot
                                         product for alignment
aligned_velocity_reward = torch.clamp(aligned_velocity, min=0.0) # Only
                                         reward positive velocities
velocity_reward = aligned_velocity_reward + velocity_penalty # Combine
                                         rewards and penalties

# Normalize velocity reward
velocity_reward = torch.clamp(velocity_reward, min=0.0)

# Combine all rewards to form total reward
total_reward = distance_reward + velocity_reward

# Components for debugging
reward_components = {
    'distance_reward': distance_reward,
    'velocity_reward': velocity_reward,
}

return total_reward, reward_components

```

718

Revolve.

```

@torch.jit.script
def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor,
left_hand_pos: torch.Tensor, right_hand_pos: torch.Tensor, left_hand_rot:
torch.Tensor, right_hand_rot: torch.Tensor) -> Tuple[torch.Tensor, Dict[str,
torch.Tensor]]:
# Define temperature variables for normalization
temp_pos = 0.1
temp_rot = 0.05
temp_catch = 0.2 # Adjusted temperature for catch reward

# Compute the distance from the object to the goal
distance_to_goal = torch.norm(object_pos - goal_pos, dim=-1)
reward_position = -distance_to_goal # Closer is better, so we take negative

# Compute the distance from the object to the left and right hands
distance_to_left_hand = torch.norm(object_pos - left_hand_pos, dim=-1) + 1e-6
# Adding a small epsilon to prevent log(0)
distance_to_right_hand = torch.norm(object_pos - right_hand_pos, dim=-1) + 1e
-6

# Encourage the object to be close to the hands with stronger feedback
reward_catch_left_base = -torch.exp(distance_to_left_hand / temp_catch) #
                                         Base feedback
reward_catch_left_prox = torch.where(distance_to_left_hand < 0.1, 1.0 -
                                         distance_to_left_hand / 0.1, torch.tensor(0.0, device=object_pos.device))
# Linear boost for proximity
reward_catch_left = reward_catch_left_base + reward_catch_left_prox # Combined rewards

reward_catch_right = -1.5 * torch.exp(distance_to_right_hand / temp_catch) #
                                         More weight for right hand penalty
reward_catch = reward_catch_left + reward_catch_right # Combined rewards for
                                         better guidance

# Compute alignment/rotation rewards based on the orientation of the hands
desired_left_rot = torch.atan2(left_hand_pos[:, 1], left_hand_pos[:, 0])
desired_right_rot = torch.atan2(right_hand_pos[:, 1], right_hand_pos[:, 0])

reward_left_rot = -torch.abs(left_hand_rot[:, 0] - desired_left_rot) #
                                         Compare some angle derived from quaternion
reward_right_rot = -torch.abs(right_hand_rot[:, 0] - desired_right_rot)

# Combine all rewards
total_reward = torch.exp(reward_position / temp_pos) + torch.exp(reward_catch
                                         / temp_catch) + torch.exp(reward_left_rot / temp_rot) + torch.exp(
                                         reward_right_rot / temp_rot)

```

719

```
# Prepare individual rewards to return
reward_components = {
    'reward_position': reward_position,
    'reward_catch': reward_catch,
    'reward_catch_left': reward_catch_left,
    'reward_catch_right': reward_catch_right,
    'reward_left_rot': reward_left_rot,
    'reward_right_rot': reward_right_rot,
}

return total_reward, reward_components
```

720

721 **NeurIPS Paper Checklist**

722 **1. Claims**

723 Question: Do the main claims made in the abstract and introduction accurately reflect the
724 paper's contributions and scope?

725 Answer: **[Yes]**

726 Justification: We accurately claim the main contributions and scope of our paper in the
727 Abstract and Introduction section.

728 Guidelines:

- 729 • The answer NA means that the abstract and introduction do not include the claims
730 made in the paper.
- 731 • The abstract and/or introduction should clearly state the claims made, including the
732 contributions made in the paper and important assumptions and limitations. A No or
733 NA answer to this question will not be perceived well by the reviewers.
- 734 • The claims made should match theoretical and experimental results, and reflect how
735 much the results can be expected to generalize to other settings.
- 736 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
737 are not attained by the paper.

738 **2. Limitations**

739 Question: Does the paper discuss the limitations of the work performed by the authors?

740 Answer: **[Yes]**

741 Justification: We discuss the limitations of our work in Sec.6.

742 Guidelines:

- 743 • The answer NA means that the paper has no limitation while the answer No means that
744 the paper has limitations, but those are not discussed in the paper.
- 745 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 746 • The paper should point out any strong assumptions and how robust the results are to
747 violations of these assumptions (e.g., independence assumptions, noiseless settings,
748 model well-specification, asymptotic approximations only holding locally). The authors
749 should reflect on how these assumptions might be violated in practice and what the
750 implications would be.
- 751 • The authors should reflect on the scope of the claims made, e.g., if the approach was
752 only tested on a few datasets or with a few runs. In general, empirical results often
753 depend on implicit assumptions, which should be articulated.
- 754 • The authors should reflect on the factors that influence the performance of the approach.
755 For example, a facial recognition algorithm may perform poorly when image resolution
756 is low or images are taken in low lighting. Or a speech-to-text system might not be
757 used reliably to provide closed captions for online lectures because it fails to handle
758 technical jargon.
- 759 • The authors should discuss the computational efficiency of the proposed algorithms
760 and how they scale with dataset size.
- 761 • If applicable, the authors should discuss possible limitations of their approach to
762 address problems of privacy and fairness.
- 763 • While the authors might fear that complete honesty about limitations might be used by
764 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
765 limitations that aren't acknowledged in the paper. The authors should use their best
766 judgment and recognize that individual actions in favor of transparency play an impor-
767 tant role in developing norms that preserve the integrity of the community. Reviewers
768 will be specifically instructed to not penalize honesty concerning limitations.

769 **3. Theory assumptions and proofs**

770 Question: For each theoretical result, does the paper provide the full set of assumptions and
771 a complete (and correct) proof?

772 Answer: **[NA]**

773 Justification: Our work does not include theoretical results.

774 Guidelines:

- 775 • The answer NA means that the paper does not include theoretical results.
- 776 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
777 referenced.
- 778 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 779 • The proofs can either appear in the main paper or the supplemental material, but if
780 they appear in the supplemental material, the authors are encouraged to provide a short
781 proof sketch to provide intuition.
- 782 • Inversely, any informal proof provided in the core of the paper should be complemented
783 by formal proofs provided in appendix or supplemental material.
- 784 • Theorems and Lemmas that the proof relies upon should be properly referenced.

785 **4. Experimental result reproducibility**

786 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
787 perimental results of the paper to the extent that it affects the main claims and/or conclusions
788 of the paper (regardless of whether the code and data are provided or not)?

789 Answer: **[Yes]**

790 Justification: We include the needed information for reproducibility in Sec.5 and Appendix
791 C.

792 Guidelines:

- 793 • The answer NA means that the paper does not include experiments.
- 794 • If the paper includes experiments, a No answer to this question will not be perceived
795 well by the reviewers: Making the paper reproducible is important, regardless of
796 whether the code and data are provided or not.
- 797 • If the contribution is a dataset and/or model, the authors should describe the steps taken
798 to make their results reproducible or verifiable.
- 799 • Depending on the contribution, reproducibility can be accomplished in various ways.
800 For example, if the contribution is a novel architecture, describing the architecture fully
801 might suffice, or if the contribution is a specific model and empirical evaluation, it may
802 be necessary to either make it possible for others to replicate the model with the same
803 dataset, or provide access to the model. In general, releasing code and data is often
804 one good way to accomplish this, but reproducibility can also be provided via detailed
805 instructions for how to replicate the results, access to a hosted model (e.g., in the case
806 of a large language model), releasing of a model checkpoint, or other means that are
807 appropriate to the research performed.
- 808 • While NeurIPS does not require releasing code, the conference does require all submis-
809 sions to provide some reasonable avenue for reproducibility, which may depend on the
810 nature of the contribution. For example
 - 811 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
812 to reproduce that algorithm.
 - 813 (b) If the contribution is primarily a new model architecture, the paper should describe
814 the architecture clearly and fully.
 - 815 (c) If the contribution is a new model (e.g., a large language model), then there should
816 either be a way to access this model for reproducing the results or a way to reproduce
817 the model (e.g., with an open-source dataset or instructions for how to construct
818 the dataset).
 - 819 (d) We recognize that reproducibility may be tricky in some cases, in which case
820 authors are welcome to describe the particular way they provide for reproducibility.
821 In the case of closed-source models, it may be that access to the model is limited in
822 some way (e.g., to registered users), but it should be possible for other researchers
823 to have some path to reproducing or verifying the results.

824 **5. Open access to data and code**

825 Question: Does the paper provide open access to the data and code, with sufficient instruc-
826 tions to faithfully reproduce the main experimental results, as described in supplemental
827 material?

828 Answer: [Yes]

829 Justification: Included with the supplementary material.

830 Guidelines:

- 831 • The answer NA means that paper does not include experiments requiring code.
- 832 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 833 • While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 834 • The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 835 • The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 836 • The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- 837 • At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 838 • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

850 6. Experimental setting/details

851 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
852 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
853 results?

854 Answer: [Yes]

855 Justification: We include these details in Sec.5, Appendix C and Appendix D.

856 Guidelines:

- 857 • The answer NA means that the paper does not include experiments.
- 858 • The experimental setting should be presented in the core of the paper to a level of detail
859 that is necessary to appreciate the results and make sense of them.
- 860 • The full details can be provided either with the code, in appendix, or as supplemental
861 material.

862 7. Experiment statistical significance

863 Question: Does the paper report error bars suitably and correctly defined or other appropriate
864 information about the statistical significance of the experiments?

865 Answer: [Yes]

866 Justification: We report the mean and standard deviation of the experiments in Tab.1 and
867 Fig. 3

868 Guidelines:

- 869 • The answer NA means that the paper does not include experiments.
- 870 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
871 dence intervals, or statistical significance tests, at least for the experiments that support
872 the main claims of the paper.
- 873 • The factors of variability that the error bars are capturing should be clearly stated (for
874 example, train/test split, initialization, random drawing of some parameter, or overall
875 run with given experimental conditions).
- 876 • The method for calculating the error bars should be explained (closed form formula,
877 call to a library function, bootstrap, etc.)
- 878 • The assumptions made should be given (e.g., Normally distributed errors).

- 879 • It should be clear whether the error bar is the standard deviation or the standard error
 880 of the mean.
 881 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 882 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 883 of Normality of errors is not verified.
 884 • For asymmetric distributions, the authors should be careful not to show in tables or
 885 figures symmetric error bars that would yield results that are out of range (e.g. negative
 886 error rates).
 887 • If error bars are reported in tables or plots, The authors should explain in the text how
 888 they were calculated and reference the corresponding figures or tables in the text.

889 **8. Experiments compute resources**

890 Question: For each experiment, does the paper provide sufficient information on the com-
 891 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 892 the experiments?

893 Answer: **[No]**

894 Justification: Partially. We report the average running time and devices of experiments
 895 description in Appendix C and E.

896 Guidelines:

- 897 • The answer NA means that the paper does not include experiments.
 898 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 899 or cloud provider, including relevant memory and storage.
 900 • The paper should provide the amount of compute required for each of the individual
 901 experimental runs as well as estimate the total compute.
 902 • The paper should disclose whether the full research project required more compute
 903 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 904 didn't make it into the paper).

905 **9. Code of ethics**

906 Question: Does the research conducted in the paper conform, in every respect, with the
 907 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

908 Answer: **[Yes]**

909 Justification: We follow the NeurIPS Code of Ethics.

910 Guidelines:

- 911 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 912 • If the authors answer No, they should explain the special circumstances that require a
 913 deviation from the Code of Ethics.
 914 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 915 eration due to laws or regulations in their jurisdiction).

916 **10. Broader impacts**

917 Question: Does the paper discuss both potential positive societal impacts and negative
 918 societal impacts of the work performed?

919 Answer: **[NA]**

920 Justification: There is no societal impact in our work.

921 Guidelines:

- 922 • The answer NA means that there is no societal impact of the work performed.
 923 • If the authors answer NA or No, they should explain why their work has no societal
 924 impact or why the paper does not address societal impact.
 925 • Examples of negative societal impacts include potential malicious or unintended uses
 926 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 927 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 928 groups), privacy considerations, and security considerations.

- 929 • The conference expects that many papers will be foundational research and not tied
930 to particular applications, let alone deployments. However, if there is a direct path to
931 any negative applications, the authors should point it out. For example, it is legitimate
932 to point out that an improvement in the quality of generative models could be used to
933 generate deepfakes for disinformation. On the other hand, it is not needed to point out
934 that a generic algorithm for optimizing neural networks could enable people to train
935 models that generate Deepfakes faster.
- 936 • The authors should consider possible harms that could arise when the technology is
937 being used as intended and functioning correctly, harms that could arise when the
938 technology is being used as intended but gives incorrect results, and harms following
939 from (intentional or unintentional) misuse of the technology.
- 940 • If there are negative societal impacts, the authors could also discuss possible mitigation
941 strategies (e.g., gated release of models, providing defenses in addition to attacks,
942 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
943 feedback over time, improving the efficiency and accessibility of ML).

944 11. Safeguards

945 Question: Does the paper describe safeguards that have been put in place for responsible
946 release of data or models that have a high risk for misuse (e.g., pretrained language models,
947 image generators, or scraped datasets)?

948 Answer: [NA]

949 Justification: There are no such risks in our paper.

950 Guidelines:

- 951 • The answer NA means that the paper poses no such risks.
- 952 • Released models that have a high risk for misuse or dual-use should be released with
953 necessary safeguards to allow for controlled use of the model, for example by requiring
954 that users adhere to usage guidelines or restrictions to access the model or implementing
955 safety filters.
- 956 • Datasets that have been scraped from the Internet could pose safety risks. The authors
957 should describe how they avoided releasing unsafe images.
- 958 • We recognize that providing effective safeguards is challenging, and many papers do
959 not require this, but we encourage authors to take this into account and make a best
960 faith effort.

961 12. Licenses for existing assets

962 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
963 the paper, properly credited and are the license and terms of use explicitly mentioned and
964 properly respected?

965 Answer: [Yes]

966 Justification: Please check in Appendix H.

967 Guidelines:

- 968 • The answer NA means that the paper does not use existing assets.
- 969 • The authors should cite the original paper that produced the code package or dataset.
- 970 • The authors should state which version of the asset is used and, if possible, include a
971 URL.
- 972 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 973 • For scraped data from a particular source (e.g., website), the copyright and terms of
974 service of that source should be provided.
- 975 • If assets are released, the license, copyright information, and terms of use in the
976 package should be provided. For popular datasets, paperswithcode.com/datasets
977 has curated licenses for some datasets. Their licensing guide can help determine the
978 license of a dataset.
- 979 • For existing datasets that are re-packaged, both the original license and the license of
980 the derived asset (if it has changed) should be provided.

- 981 • If this information is not available online, the authors are encouraged to reach out to
982 the asset's creators.

983 **13. New assets**

984 Question: Are new assets introduced in the paper well documented and is the documentation
985 provided alongside the assets?

986 Answer: **[Yes]**

987 Justification: Included with the supplementary material.

988 Guidelines:

- 989 • The answer NA means that the paper does not release new assets.
990 • Researchers should communicate the details of the dataset/code/model as part of their
991 submissions via structured templates. This includes details about training, license,
992 limitations, etc.
993 • The paper should discuss whether and how consent was obtained from people whose
994 asset is used.
995 • At submission time, remember to anonymize your assets (if applicable). You can either
996 create an anonymized URL or include an anonymized zip file.

997 **14. Crowdsourcing and research with human subjects**

998 Question: For crowdsourcing experiments and research with human subjects, does the paper
999 include the full text of instructions given to participants and screenshots, if applicable, as
1000 well as details about compensation (if any)?

1001 Answer: **[NA]**

1002 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1003 Guidelines:

- 1004 • The answer NA means that the paper does not involve crowdsourcing nor research with
1005 human subjects.
1006 • Including this information in the supplemental material is fine, but if the main contribu-
1007 tion of the paper involves human subjects, then as much detail as possible should be
1008 included in the main paper.
1009 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1010 or other labor should be paid at least the minimum wage in the country of the data
1011 collector.

1012 **15. Institutional review board (IRB) approvals or equivalent for research with human
1013 subjects**

1014 Question: Does the paper describe potential risks incurred by study participants, whether
1015 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1016 approvals (or an equivalent approval/review based on the requirements of your country or
1017 institution) were obtained?

1018 Answer: **[NA]**

1019 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1020 Guidelines:

- 1021 • The answer NA means that the paper does not involve crowdsourcing nor research with
1022 human subjects.
1023 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1024 may be required for any human subjects research. If you obtained IRB approval, you
1025 should clearly state this in the paper.
1026 • We recognize that the procedures for this may vary significantly between institutions
1027 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1028 guidelines for their institution.
1029 • For initial submissions, do not include any information that would break anonymity (if
1030 applicable), such as the institution conducting the review.

1031 **16. Declaration of LLM usage**

1032 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1033 non-standard component of the core methods in this research? Note that if the LLM is used
1034 only for writing, editing, or formatting purposes and does not impact the core methodology,
1035 scientific rigorousness, or originality of the research, declaration is not required.

1036 Answer: [Yes]

1037 Justification: Please refer to Sec.4 and Appendix B for detailed usage.

1038 Guidelines:

- 1039 • The answer NA means that the core method development in this research does not
1040 involve LLMs as any important, original, or non-standard components.
- 1041 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1042 for what should or should not be described.