

AUTOMATIC REWARD SHAPING FROM MULTI-OBJECTIVE HUMAN HEURISTICS

Anonymous authors

Paper under double-blind review

ABSTRACT

Designing effective reward functions remains a central challenge in reinforcement learning, especially in multi-objective environments. In this work, we propose **Multi-Objective Reward Shaping with Exploration**, a general framework that automatically combines multiple human-designed heuristic rewards into a unified reward function. MORSE formulates the shaping process as a bi-level optimization problem: the inner loop trains a policy to maximize the current shaped reward, while the outer loop updates the reward function to optimize task performance. To encourage exploration in the reward space and avoid suboptimal local minima, MORSE introduces stochasticity into the shaping process, injecting noise guided by task performance and the prediction error of a fixed, randomly initialized neural network. Experimental results in MuJoCo and Isaac Sim environments show that MORSE effectively balances multiple objectives across various robotic tasks, achieving task performance comparable to those obtained with manually tuned reward functions.

1 INTRODUCTION

Reward design remains a fundamental challenge in robot learning, particularly when multiple conflicting objectives must be balanced. A typical task definition might involve the robot traversing a specified distance (e.g., 2 meters) within a given time limit (e.g., 10 seconds) while maintaining the torque below a safety limit (e.g., 45 Nm). While the success criteria for the task are straightforward to define, their sparsity makes them inadequate as rewards for effective policy optimization. (Andrychowicz et al., 2018)

To facilitate policy learning, researchers often design dense heuristic reward functions, such as rewarding higher velocities or penalizing large joint movements to minimize energy cost. However, these heuristics frequently conflict with one another, necessitating labor-intensive manual tuning to combine them into a final reward function. Studies show that over 90% of reinforcement learning (RL) practitioners devote substantial effort to adjusting reward functions, highlighting a critical scalability bottleneck (Booth et al., 2023).

Our work proposes a self-supervised reward shaping framework that automates the reward design process. The approach requires RL practitioners to specify only: (1) a task performance criterion, which scores a trajectory at the end of the episode, and (2) a set of heuristic reward functions, which provide dense feedback for each step. The system then automatically learns an optimal combination of these heuristics through a bi-level optimization process (Zhang et al., 2024b), where the inner loop trains RL policies under the current reward function, and the outer loop adapts the reward function itself to maximize task performance.

However, it is not trivial to apply bi-level optimization to shape the rewards for real-world robotic tasks. Consider the robotic dog benchmark in (Margolis & Agrawal, 2023), which defines 15 distinct objectives, constituting an enormous space of valid reward combinations. With such a large reward weight space, the weight-performance landscape becomes highly non-convex with numerous local extrema (Xie et al., 2024). Under such circumstances, existing bi-level methods (Gupta et al., 2023) fail to converge reliably in complex scenarios, as we will further demonstrate in the experiments.

To overcome these challenges, we introduce **Multi-Objective Reward Shaping with Exploration** (MORSE). MORSE augments the outer loop of bi-level optimization with controlled stochastic ex-

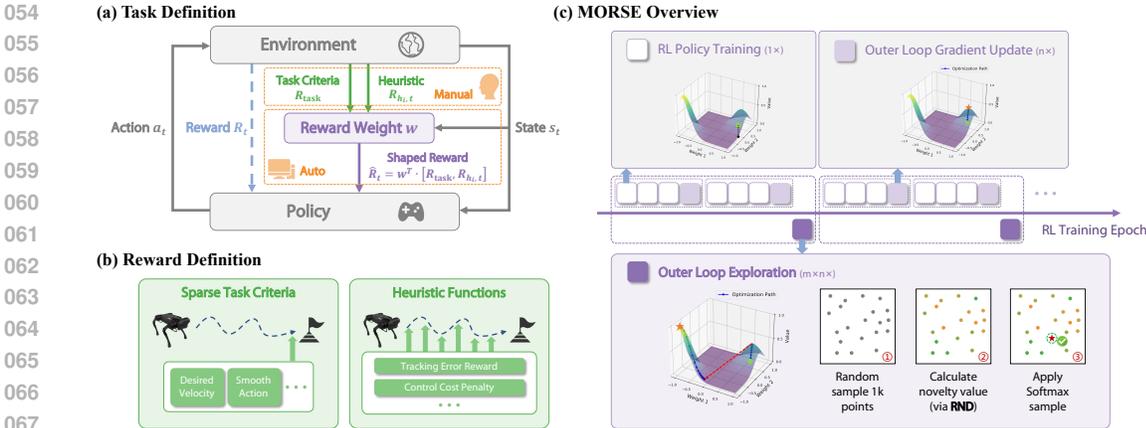


Figure 1: (a) Experts provide only task criteria and heuristic functions, instead of a manually tuned reward function. (b) Example task criteria and heuristic functions for robotic dog. (c) The inner loop trains the RL policy on shaped rewards, while the outer loop, combining gradient updates with exploration, updates reward weights to maximize task criteria. In the 3D graphs, we use xy plane to denote reward weight space, and z axis to represent the corresponding task performance.

ploration, thereby facilitating training when the reward space is large. Concretely, if the policy plateaus without reaching the target objective, the algorithm initiates outer-loop exploration: it samples a new reward weight and resumes gradient-based bi-level updates from this point. Exploration relies on novelty scores computed via Random Network Distillation (RND)(Burda et al., 2018), where the prediction error against a fixed, randomly initialized network serves as the novelty measure. New starting points are then selected through softmax sampling over these scores.

The experiments are organized as follows. We begin with a motivating example in a custom multi-objective CartPole environment and show that unguided exploration helps conventional bi-level approaches to escape from local optima. We then evaluate MORSE in two complementary settings. (1) Synthetic optimization surfaces, where we simplify bi-level optimization into single-level gradient update to isolate and validate each component of the controlled exploration mechanism. (2) Challenging robotics domains in MuJoCo (Towers et al., 2024) and Isaac Lab(Mittal et al., 2023), including a 9-objective quadruped locomotion task and a 5-objective robotic manipulation task. Across all environments, MORSE consistently outperforms vanilla bi-level optimization and achieves performance comparable to that obtained with human-engineered reward functions.

2 RELATED WORK

2.1 MULTI-OBJECTIVE REINFORCEMENT LEARNING

To incorporate the trade-offs between multiple objectives, different multi-objective reinforcement learning (MORL) methods have been proposed. These methods can be roughly divided into two categories, that is, single-policy and multi-policy (Hayes et al., 2022; Felten et al., 2023).

Single-policy methods yield a policy capable of adapting to any preference utility at test time. This can be done either during training, where the policy maximizes expected scalarized returns (ESR) for potential utilities, (Roijers et al., 2018; Siddique et al., 2020) or during execution, where the policy selects actions based on test-time utility(Van Moffaert et al., 2013a;b; Issabekov & Vamplew, 2012; Reymond et al., 2023).

Multi-policy methods, on the other hand, generate a collection of policies that form a Pareto Front or Coverage Set, capturing inter-objective relationships. They can iteratively sample utility functions and derive corresponding policies to construct an exhaustive coverage set (Parisi et al., 2014; Mossalam et al., 2016; Parisi et al., 2017; Natarajan & Tadepalli, 2005), or generate multiple policies concurrently and employ specific criteria such as hypervolume value or Monte Carlo Tree Search to select actions(Castelletti et al., 2012; Ruiz-Montiel et al., 2017; Van Moffaert & Nowé, 2014; Yang et al., 2019; Lu et al., 2022; Reymond et al., 2022).

In general, MORL focuses on recovering the Pareto front that captures the trade-offs between multiple objectives, but it does not learn how to choose or optimize the reward weights. It assumes that the appropriate weights are given at test time. Our method aims to solve the orthogonal problem. Our goal is to automatically adjust the reward weights of auxiliary rewards to better guide the RL policy toward higher success rates. Therefore, our problem setup fundamentally differs from MORL.

2.2 REWARD SHAPING IN REINFORCEMENT LEARNING

Designing effective reward functions poses fundamental challenges in RL: sparse rewards yield inefficient learning, while poorly shaped rewards may induce suboptimal behaviors. To address these challenges, existing reward shaping methods can be broadly categorized into four main approaches.

Heuristic-based intrinsic reward methods encourage exploration through visitation frequency counts (Tang et al., 2017; Lobel et al., 2023; Machado et al., 2020; Fu et al., 2017) or curiosity-driven uncertainty bonuses (Pathak et al., 2017; 2019; Zhang et al., 2020; Burda et al., 2018), but require careful manual tuning of intrinsic and extrinsic reward balances. More automated approaches include differentiable intrinsic reward learning (LIRPG Zheng et al. (2018)), trajectory preference-based reward inference (SORS Memarian et al. (2021)), and composite reward functions combining task and exploration objectives (EXPLORS Devidze et al. (2022)). With the rapid development of large language models (LLMs), several recent works have explored using LLMs to assist in constructing reward functions (Ma et al., 2023; Zhang et al., 2024a). These approaches typically require the LLM to define both the heuristic functions and the reward weights, and to iteratively refine them over multiple optimization rounds, which leads to substantial computational overhead.

Most relevant to our work are bi-level optimization methods (Hu et al., 2020; Gupta et al., 2023) that formulate reward shaping as $\max_{\phi} \mathcal{R}_{task}(\pi_{\theta}^*)$ s.t. $\pi_{\theta}^* = \arg \max_{\pi} \mathcal{R}_{\phi}(s, a)$, where ϕ parameterizes the shaped reward. These methods assume (i) smooth reward-performance landscapes and (ii) sufficient gradient signals from heuristic rewards. However, such assumptions are often violated in robotics tasks, causing outer-loop optimization to stagnate in local optima. Our method addresses these limitations by adding stochasticity that considers task performance and novelty measurement in reward weight space.

3 PRELIMINARY

3.1 TASK FORMULATION

We formulate each task as a Markov decision process (MDP) defined by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \hat{R}, \gamma)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, \mathcal{P} is the state transition function, $\hat{R}: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the task-specific reward function, γ is the discount factor. The goal is to find a policy π_{θ} that maximizes the discounted cumulative task reward, which is defined as $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \gamma^t \hat{R}(s_t, a_t)]$, where $\tau = (s_0, a_0, \dots, s_T)$ denotes a trajectory sampled under π_{θ} , and T is the episode horizon.

Designing an effective reward function \hat{R} is challenging, as it often requires balancing multiple conflicting objectives. In the proposed framework, we gradually learn \hat{R} from (1) a sparse task-completion metric R_{task} provided by a human expert (e.g., a binary indicator of success at episode termination) and (2) a set of n auxiliary heuristic reward functions $\{R_{h_i}\}_{i=1}^n$, each capturing a separate objective of the task. Here, the human expert only needs to determine the terminal success criterion and provide heuristic guidance; the algorithm automatically finishes reward shaping.

Formally, the algorithm balances all reward components $\mathbf{R} = (R_{task}, R_{h_1}, \dots, R_{h_n})$ to maximize task reward R_{task} . We learn a state-dependent weight function $w_{\phi}: \mathcal{S} \rightarrow \mathbb{R}^{n+1}$ to produce the composite reward:

$$R_{\phi}(s, a) = w_{\phi}(s)^{\top} \mathbf{R}(s, a),$$

where $\mathbf{R}(s, a) \in \mathbb{R}^{n+1}$ is the concatenated vector of all reward components.

3.2 BI-LEVEL OPTIMIZATION FOR REINFORCEMENT LEARNING

The above reward shaping process can be solved via bi-level optimization. The outer loop optimizes the reward weight $w_\phi(s)$ to maximize the accumulated task reward $J(\theta, \phi)$:

$$J(\theta, \phi) = E_{\pi_\theta, w_\phi} [\sum_{t=0}^T \gamma^t R_{\text{task}}(s_t, a_t)]$$

where as the inner loop optimizes the policy π_θ to maximize the weighted return $G(\theta)$:

$$G_\phi(\theta) = E_{\pi_\theta} [\sum_{t=0}^T \gamma^t R_\phi(s, a)] = E_{\pi_\theta} [\sum_{t=0}^T \gamma^t w_\phi(s_t) \cdot \mathbf{R}(s_t, a_t)]$$

The inner loop can be optimized using any reinforcement learning algorithm, while the outer loop can be solved using various bi-level optimization techniques. In this work, we choose the Implicit Function (IF)-based approach (Zhang et al., 2024b).

According to the chain rule, the gradient computation for the outer loop can be decomposed as:

$$\frac{dJ(\theta, \phi)}{d\phi} = \frac{\partial J(\theta, \phi)}{\partial \phi} + \frac{\partial J(\theta, \phi)}{\partial \theta} \cdot \frac{d\theta}{d\phi}$$

Since w_ϕ does not explicitly influences $J(\theta, \phi)$, the term $\frac{\partial J(\theta, \phi)}{\partial \phi}$ equals zero. The next term $\frac{\partial J(\theta, \phi)}{\partial \theta}$ is the policy gradient calculated using the task reward, R_{task} . The third term $\frac{d\theta}{d\phi}$ captures how the optimal inner loop policy responds to changes in the reward weight w_ϕ . It is determined by $\frac{d\pi_{\theta; \phi}^*}{d\phi}$, where $\pi_{\theta; \phi}^*$ is the optimal policy obtained in the inner loop with w_ϕ .

We briefly discuss how to calculate the second term. During convergence, we have $\nabla_\theta G_\phi(\theta) = 0$, leading to the following expression for $\frac{d\pi_{\theta; \phi}^*}{d\phi}$:

$$\frac{d\pi_{\theta; \phi}^*}{d\phi} = - \left(\nabla_\phi^2 \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t R_\theta \right] \right)^{-1} \nabla_\phi \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t R_\theta \right]$$

The two terms above can be computed using the second-order and first-order derivatives of the inner loop objective function. However, computing the inverse Hessian matrix $H = \nabla_\phi^2 \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t R_\theta \right]$ is computationally expensive. Therefore, we use the Neumann Series method to approximate the inverse of the Hessian matrix :

$$H^{-1} \approx \sum_{i=0}^{K-1} (I - H)^i$$

4 A MOTIVATION EXAMPLE: MULTI-OBJECTIVE CARTPOLE WITH EXTREMELY SPARSE REWARD

In this section, we extend the classic CartPole to a multi-objective setting where the agent must balance the pole while moving the cart as far right as possible. This demonstrates how vanilla bi-level optimization struggles with conflicting objectives and highly sparse task rewards. The environment retains the original observation and action spaces but uses a composite reward with four components: (1) a +100 task reward for surviving all 100 steps; (2) a +1 survival heuristic per step; (3) a +1 position heuristic when the cart’s position exceeds 0.5; and (4) a -1 interference penalty when the agent’s action matches a PD controller’s stabilization attempt. These modifications make the multi-objective Cartpole significantly harder. We perform a grid search for potential reward weights, and the best performing configuration achieved approximately a 50% success rate.

We use REINFORCE in the inner loop for policy training, with the outer loop modeling rewards as $w(s) \cdot \mathbf{R}(s, a)$ through a weight network $w(s)$ that maps states to reward weights. We compare three update strategies for $w(s)$: **Constant** (randomly initialized weights with no update), **Gradient** (updated with standard bi-level optimization), and **Gradient w/ Reset** (resetting the weight network parameters every five gradient updates).

We evaluate each method over 10 trials with different random seeds, and report the results in Fig. 2a. Compared with **Constant**, **Gradient** improves the task success rate to 30%, while **Gradient w/ Reset**, which periodically resets the reward weight, further increases the success rate to 50%.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

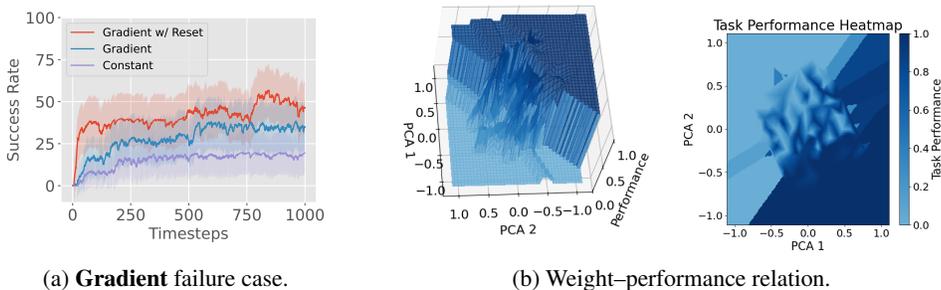


Figure 2: Lesson from the motivation example: updating reward weights with outer-loop gradient is not sufficient; periodic resetting reward weights helps to escape from poor local optima.

The limitation of **Gradient** arises from two inherent properties of multi-objective problems: (1) adding competing objectives increases task complexity and expands the search space for reward weights; and (2) task criteria provide extremely sparse feedback, yielding weak gradient signals for backpropagation. Starting with a poor initial weight often leads the policy to converge to a sub-optimal solution, which in turn prevents the outer loop from making meaningful gradient updates. Consequently, gradient-based optimization is only effective when the number of objectives is small, the reward space is limited, the solution space is broad, and the mapping from weights to performance is smooth and continuous, avoiding rugged, non-convex landscapes.

However, real-world tasks rarely satisfy these conditions. To illustrate this, we visualize the Walker2d task in the MuJoCo domain in Fig. 2b. Since Walker2d has three objectives, we project the 3D reward weight space into 2D via Principal Component Analysis (PCA), and use the z-axis to represent the optimal task performance achieved under each weight. We randomly sample around 100 weights in the scenario, use these weights to train separate RL policies, and record the final performance of each policy. After interpolating the task performance, we construct the performance surface. The resulting landscape exhibits numerous local optima, violating the assumptions under which **Gradient** succeeds.

Interestingly, our experiments show that randomly resetting reward weights helps escape poor local optima, leading to substantially better performance. We will provide more rigorous empirical validation of these findings in Sec. 6.1.

5 MORSE: MULTI-OBJECTIVE REWARD SHAPING WITH EXPLORATION

The motivation example suggests that introducing noise into the bi-level optimization can potentially discover superior reward function and improve convergence performance. Inspired by this insight, we present a general framework to enhance bi-level optimization and propose a practical technique, MORSE. In this section, we present MORSE, a stochastic optimization method that enhances exploration in the outer loop through selective noise injection based on two metrics, reward space novelty and task performance.

5.1 EXPLORATION GUIDED BY TASK PERFORMANCE AND REWARD SPACE NOVELTY

We present the outer-loop exploration mechanism in Alg. 1.

Each episode receives a task reward of 1 if it satisfies task success criteria, otherwise 0. We use P to denote task performance, which is the fraction of successful trajectories over all rollouts obtained using the current policy. Exploration is triggered with probability $1 - P$.

Reward-space novelty is measured using a curiosity-driven metric based on Random Network Distillation (Burda et al., 2018). Specifically, we introduce two networks: a fixed, randomly initialized target network f_{target} and a trainable predictor f_{pred} , both mapping reward weights to scalar values. Novelty is quantified as the prediction error between them, with larger discrepancies indicating more novel reward weights. The predictor is trained with the reward weights previously used for RL

training, minimizing the MSE loss $= \|f_{\text{pred}}(w(s)) - f_{\text{target}}(w(s))\|_2$. We define the novelty metric for a weight candidate w as $V_{\text{novelty}}(w) = \|f_{\text{pred}}(w(s)) - f_{\text{target}}(w(s))\|_2$. Larger score indicates that w lies in a less explored region.

Outer-loop exploration proceeds by randomly sampling N candidate reward weights $W = \{w_i\}_{i=1}^N$ from the weight space. Each candidate is evaluated by the novelty metric $V_{\text{novelty}}(w_i)$, and a new weight W' is selected via softmax sampling over these scores.

Algorithm 1 Outer Loop Exploration

Require: History reward weights W_{past} , sample size N , the current task performance R_{curr} , the variation in task performance between two training epochs ΔR , improvement threshold α

- 1: **if** $\Delta R < \alpha$ **then**
 - 2: Train the predictor network f_{pred} with $\|f_{\text{pred}}(w) - f_{\text{target}}(w)\|_2, w \in W_{\text{past}}$
 - 3: Let $P_{\text{curr}} = \text{normalize}(R_{\text{curr}})$. With probability $1 - P_{\text{curr}}$, perform the following steps:
 - 4: Randomly sample N points $W = \{w_i\}_{i=1}^N$ from the reward weight space
 - 5: Compute novelty metric: $V_{\text{novelty}}(w_i), \forall w_i \in W$
 - 6: Sample from W based on softmax probabilities of V_{novelty} : $w \sim \text{Softmax}(V_{\text{novelty}})$
 - 7: **end if**
-

5.2 OVERALL ALGORITHM

Combining the bi-level optimization and outer loop exploration, we present a realization of Multi-Objective Reward Shaping with Exploration, which is summarized in Alg. 2.

In each training epoch, we rollout data under the current policy π_θ and update θ via any on-policy or off-policy RL optimizer. Every T_{grad} epochs, the multi-objective weight vector ϕ is updated by ascending the gradient of $J(\phi)$, and, at a lower explore frequency T_{explore} , we reset ϕ through outer loop exploration and reinitialize the policy’s parameters to encourage the policy to retain high entropy. By jointly optimizing θ and ϕ while periodically reinitializing the policy, MORSE balances efficient exploitation of learned reward compositions with continual exploration in the multi-objective space, yielding robust policies that negotiate conflicting objectives without manual reward engineering.

Algorithm 2 Multi-Objective Reward Shaping with Exploration (MORSE)

Require: Total training epoch N_{epoch} , RL policy π_θ , outer loop update interval T_{grad} , exploration interval T_{explore}

- 1: **for** $e = 0$ to N_{epoch} **do**
 - 2: **Rollout:** Execute policy π_θ to collect experience
 - 3: **Update policy:** Perform RL updates to optimize π_θ
 - 4: **if** $e \bmod T_{\text{grad}} = 0$ **then**
 - 5: Update reward weight: $\phi_{e+1} \leftarrow \phi_e + \nabla_\phi J(\phi_e)$
 - 6: **if** $e \bmod T_{\text{explore}} = 0$ **then**
 - 7: $w_\phi \leftarrow \text{Outer_Loop_Exploration}(\phi_e)$
 - 8: **Reset policy:** Reset the whole actor network parameters and the last layer of critic network
 - 9: **end if**
 - 10: **end if**
 - 11: **end for**
-

6 EXPERIMENTS

6.1 PRELIMINARY VALIDATION OF OUTER LOOP EXPLORATION

We validate the effect of outer loop in a simplified environment. Focusing only on the outer loop update and ignoring the inner loop RL training, we treat reward shaping as an optimization problem in a multi-dimensional space. Each reward weight corresponds to a point in the space, and the optimization objective is task performance, represented as a scalar value. For simplicity, we assume

the reward weight is $2D$, $x \in [-1, 1]^2$, and task performance is determined by the function $f(x) \in [0, 1]$. The optimization goal is to find the input value that maximizes $f(x)$ within 100 steps.

We design three types of functions with different characteristics to construct varying optimization landscapes: a smooth polynomial function with multiple local optima (SmoothPolynomial), a complex fixed neural network function (FixedNN), and a function with sparse spike regions (RandomSpiky). These functions are normalized to ensure the output is always within the $[0, 1]$ range. For each of the three function types, we randomly generate 10 parameter sets, resulting in a total of 30 functions. We conduct 10 independent optimization runs for each function, starting from different random initial points.

6.1.1 TIMING OF EXPLORATION

We evaluate several strategies for resetting the reward weight. **No Exploration** updates the reward weights only by the gradient of the outer loop with no exploration. **Periodic Exploration** samples new reward weights based on novelty scores and resets the reward weights after every 10 gradient updates. **MORSE** reset the reward weight with the novelty metric and performance criteria after every 10 gradient updates.

	RandomSpiky	SmoothPolynomial	FixedNN	Average
No Exploration	0.200 (0.422)	0.832 (0.213)	0.945 (0.141)	0.659 (0.432)
Periodic Exploration	0.102 (0.316)	0.690 (0.214)	0.695 (0.306)	0.496 (0.393)
MORSE	0.873 (0.319)	0.799 (0.273)	0.946 (0.140)	0.873 (0.254)

Table 1: Average performance of different exploration timing across 10 seeds. Each cell reports the mean performance, with standard deviation shown in parentheses. **MORSE** outperforms other baselines in two of the three optimization landscapes.

As shown in Tab. 1, the results show that **MORSE** significantly outperforms all baseline methods. **No Exploration** converges to suboptimal solutions on complex optimization landscapes, such as RandomSpiky. Meanwhile, **Periodic Exploration** reset the reward weights periodically, disrupting the gradient-based optimization process, which leads to worse performance compared to the simple gradient-based method, **No Exploration**. In contrast, **MORSE** can terminate random exploration based on performance criteria, achieving a balance between exploration and gradient optimization.

6.1.2 SAMPLING ALGORITHM FOR EXPLORATION

In this set of experiments, we replace RND-based sampling with alternative strategies while keeping all other settings fixed. The baselines include random sampling (random) and two canonical sample-based optimization methods, the cross-entropy method (CEM) (Rubinstein, 1999) and covariance matrix adaptation (CMA) (Hansen et al., 2003). For CEM and CMA, each batch samples five candidate weight vectors, and the top 40% (two elites) are retained to guide the next round.

	RandomSpiky	SmoothPolynomial	FixedNN	Average
random	0.701 (0.483)	0.872 (0.204)	0.946 (0.140)	0.840 (0.320)
CEM	0.302 (0.482)	0.711 (0.273)	0.770 (0.186)	0.594 (0.388)
CMA	0.801 (0.422)	0.698 (0.329)	0.932 (0.185)	0.810 (0.330)
RND(ours)	0.873 (0.319)	0.799 (0.273)	0.946 (0.140)	0.873 (0.254)

Table 2: Average performance of different sampling strategies across 10 seeds. Each cell reports the mean performance, with standard deviation shown in parentheses. RND-based exploration outperforms other baselines in two of the three optimization landscapes.

We present the results in Tab. 2. Compared with random, RND discovers more novel reward weights within a limited number of exploration steps. The performance of CEM and CMA degrades significantly when the number of samples is restricted, as in our setting, or when the optimization landscape is highly non-smooth, as in RandomSpiky. Both methods generate multiple candidates per round and update their sampling strategies based on batch performance. In contrast, our approach samples weights sequentially, one at a time, which substantially reduces the total number of samples required to reach a feasible solution. Since each sample corresponds to a full RL training process, this results in markedly higher computational efficiency.

6.2 MORSE IN RL SCENARIOS

We evaluate MORSE in 8 robotic tasks: 3 MuJoCo locomotion tasks, *Halfcheetah*, *Hopper*, and *Walker2d*, and 5 Isaac Sim tasks, including 3 locomotion tasks, *Quadcopter* (Isaac-Quadcopter-Direct-v0), *Unitree-A1-2obj* (modified Isaac-Velocity-Flat-Unitree-A1-v0), *Unitree-A1-9obj* (original Isaac-Velocity-Flat-Unitree-A1-v0), and 2 manipulation tasks, *Reach* (Isaac-Reach-Franka-v0), *Lift* (Isaac-Lift-Cube-Franka-IK-Rel-v0). These tasks involve 2, 3, 3, 3, 2, 9, 4, and 5 reward objectives, respectively.

6.2.1 EXPERIMENT SETUP

Task Criteria and Heuristic Functions. An episode is considered successful only if the agent survives to the end, reaches the target goal (e.g., specific velocity or position), and its entire trajectory remains within a predefined action threshold. In our experiments, the action threshold is determined by training a policy with oracle weights defined in the codebase and measuring the maximum action features (e.g., torque or angular velocity) of the policy rollouts. For practical applications, these action thresholds are often set empirically, as in (Chen et al., 2025). Heuristic rewards correspond to the single-objective reward components defined in the original codebase. Detailed definitions of task criteria and heuristic functions for each environment are provided in the Appendix.

Implementation Details. Given the distinct API and rollout logic between MuJoCo and Isaac Sim, we implement MORSE using two different frameworks, Stable-Baselines3(Raffin et al., 2021) for MuJoCo and rsl-rl(Rudin et al., 2022) for Isaac Sim. In Stable-Baselines3, PPO’s critic is extended to a multi-head critic that predicts returns for each heuristic function. In rsl-rl, we keep a single-head critic but store all reward components in the buffer.

Reward Weight Constraints. For MuJoCo tasks, we design two difficulty levels. **Easy** setting initialize reward weights randomly within $[0, 1]$ and maintains this range for weight update. In **Hard**, the range of reward weight is $[-1, 1]$. The larger range makes it significantly more challenging to identify the optimal reward function. For Isaac Sim tasks, weights are initialized within $[0, 1]$ and clipped accordingly.

All experiments are repeated 8 times with different random seeds. In the plots, solid lines in the plots denote the means, and shaded regions denote the standard deviations. For more details on hyperparameter settings, please refer to the Appendix.

6.2.2 MAIN RESULTS

We compare the following methods. **Oracle:** The policy is trained with the constant reward weight tuned by human expert. **Gradient:** The reward weight is updated using only the bi-level optimization gradient, equivalent to Barfi(Gupta et al., 2023). **MORSE:** The reward weights are reset with MORSE.

Fig. 3 reports the main experiment results. Across all environments, **MORSE** consistently outperforms **Gradient** and achieves performance comparable to **Oracle**, successfully identifying reward weights that enable high policy success rates. Furthermore, the results on MuJoCo-easy and MuJoCo-hard tasks show that, as the reward space expands, the performance of **Gradient** degrades sharply, whereas **MORSE** maintains stable performance. This again validates the insights in Sec. 4: gradient-based optimization method is strongly sensitive to the initial weight selection and the size of exploration space, and is also prone to suboptimal convergence, whereas exploration helps it escape local optima. These findings highlight the critical role of outer-loop exploration in complex robotic tasks and demonstrate the potential of **MORSE** as a general framework for reward shaping.

6.2.3 ABLATION STUDIES

MORSE incorporates three key design choices: (1) during gradient-based bi-level optimization, if the policy converges to unsatisfactory performance, the algorithm jumps to a new reward weight; (2) during exploration, the new starting weight is sampled based on the RND novelty metric; (3) when resetting the reward weight, the actor of the RL policy is also reset to encourage sufficient policy exploration.

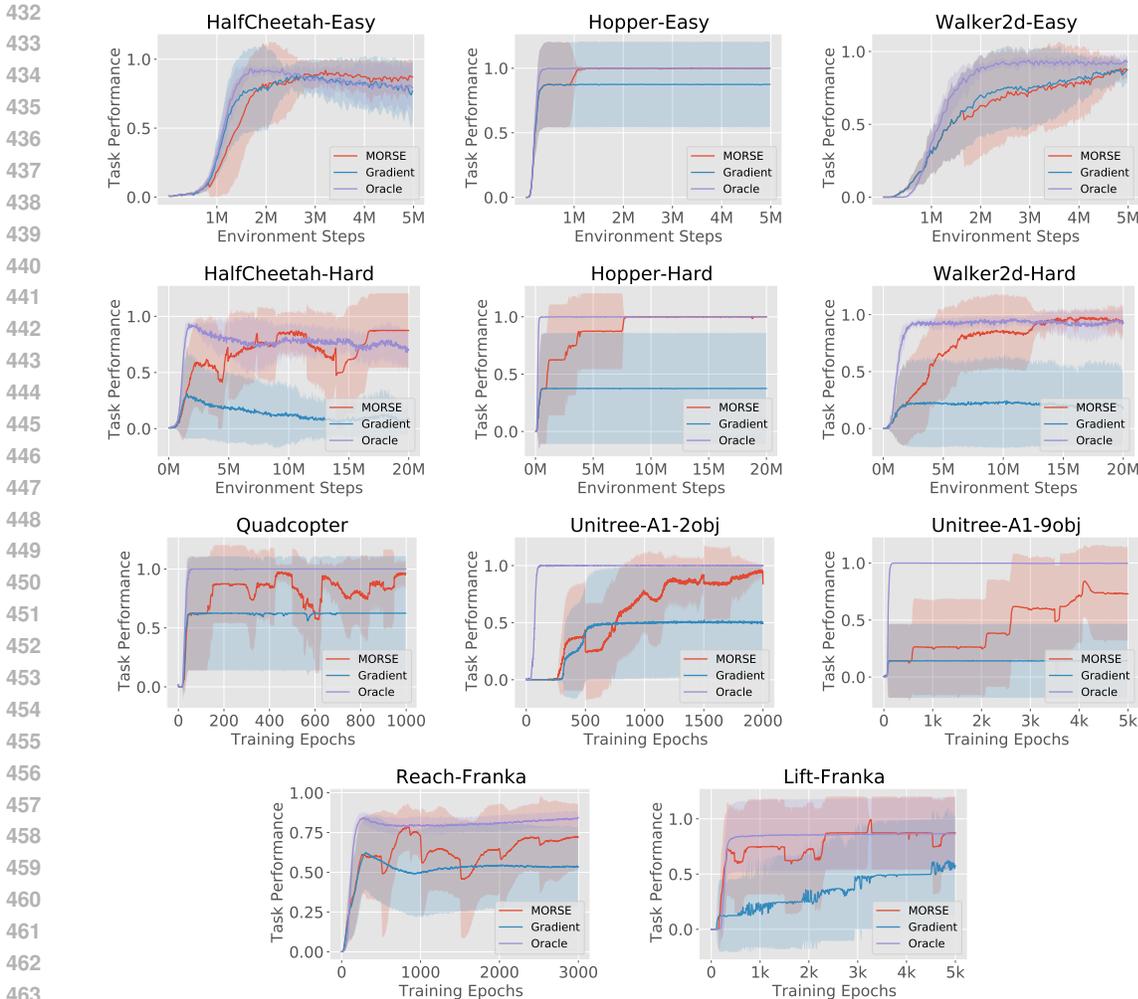


Figure 3: Main results in the MuJoCo and Isaac Sim environments. The results show that **MORSE** outperforms the vanilla bi-level optimization method and is comparable with **Oracle**, whose reward weights are tuned by human expert.

We conduct an ablation study in the MuJoCo-Hard domain, comparing **MORSE** with four variants: **wo/ Gradient**, which conducts outer loop exploration but does not perform gradient update; **wo/ Performance Metric**, which resets weights at fixed intervals, regardless of policy performance; **wo/ Novelty Metric**, which replaces RND with random sampling during exploration; and **wo/ Reset Policy**, which keeps the inner loop policy unchanged when resetting reward weights. Results are shown in Fig. 4.

wo/ Gradient performs worse than **MORSE**, suggesting that outer-loop gradient updates help with weight selection. However, **Gradient Only** suggests that, gradient updates are insufficient and must be complemented with outer loop exploration mechanism.

The discrepancy between **MORSE** and **wo/ Performance Metric** highlights the importance of reset conditioning based on task performance. In **wo/ Performance Metric**, the reward weight switches periodically even after it reaches the global optimum, which degrades final performance.

wo/ Novelty Metric randomly samples a new reward weight during exploration. For simpler tasks (Hopper) where many reward weights yield high performance, the novelty metric has less impact. However, for harder tasks like Walker2D, selecting the most novel reward weight speeds up the coverage of the whole reward space.

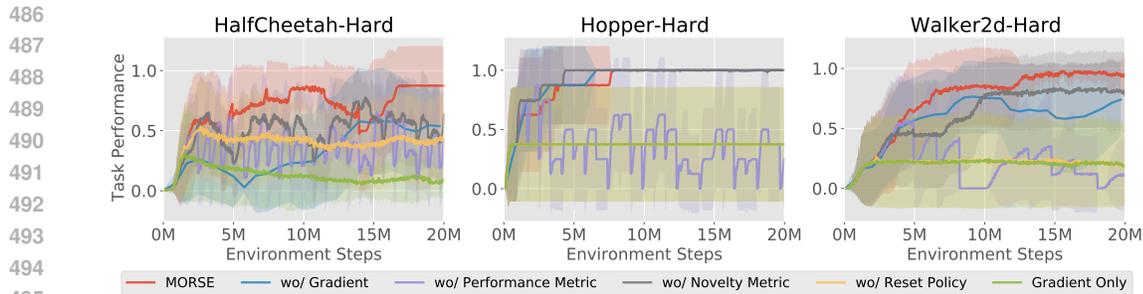


Figure 4: Ablate key components of MORSE in MuJoCo hard domain.

Comparing **wo/ Reset Policy** with **MORSE**, we find that resetting the policy during outer loop exploration is essential for the policy to maintain high entropy and therefore adapt to reward changes.

7 CONCLUSION AND FUTURE WORK

Our goal is to reduce the burden of manual reward shaping in RL, simplifying both task design and training. In MORSE, human experts only need to specify sparse task criterion and dense heuristic functions without defining their interactions, while the policy automatically learns to solve complex robotic tasks. We hope this framework can make RL training more scalable and broadly applicable.

There remain several directions for improvement. In outer-loop exploration, novelty and performance metrics could be combined more effectively, for example by considering the historical performance of reward weights and adopting Upper Confidence Bound (Auer et al., 2002) to better balance exploration and exploitation. Moreover, the current design requires resetting the whole policy when switching reward weights. Incorporating successor features (Barreto et al., 2017) into the policy structure may enable more efficient reuse of past knowledge and faster adaptation to new reward weights.

8 REPRODUCIBILITY STATEMENT

We have uploaded the source codes in the supplemental materials to ensure reproducibility.

REFERENCES

- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay, 2018. URL <https://arxiv.org/abs/1707.01495>.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- Serena Booth, W Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi. The perils of trial-and-error reward design: Misdesign through overfitting and invalid task specifications. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence (AAAI)*, Feb 2023.
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- Andrea Castelletti, Francesca Pianosi, and Marcello Restelli. Tree-based fitted q-iteration for multi-objective markov decision problems. In *The 2012 international joint conference on neural networks (IJCNN)*, pp. 1–8. IEEE, 2012.

- 540 Jiayu Chen, Chao Yu, Yuqing Xie, Feng Gao, Yinuo Chen, Shu'ang Yu, Wenhao Tang, Shilong Ji,
541 Mo Mu, Yi Wu, et al. What matters in learning a zero-shot sim-to-real rl policy for quadrotor
542 control? a comprehensive study. *IEEE Robotics and Automation Letters*, 2025.
- 543 Rati Devidze, Parameswaran Kamalaruban, and Adish Singla. Exploration-guided reward shaping
544 for reinforcement learning under sparse rewards. *Advances in Neural Information Processing*
545 *Systems*, 35:5829–5842, 2022.
- 546 Florian Felten, Lucas N. Alegre, Ann Nowé, Ana L. C. Bazzan, El Ghazali Talbi, Grégoire Danoy,
547 and Bruno Castro da Silva. A toolkit for reliable benchmarking and research in multi-objective
548 reinforcement learning. In *Proceedings of the 37th Conference on Neural Information Processing*
549 *Systems (NeurIPS 2023)*, 2023.
- 550 Justin Fu, John Co-Reyes, and Sergey Levine. Ex2: Exploration with exemplar models for deep
551 reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- 552 Dhawal Gupta, Yash Chandak, Scott Jordan, Philip S Thomas, and Bruno C da Silva. Behavior
553 alignment via reward function optimization. *Advances in Neural Information Processing Systems*,
554 36:52759–52791, 2023.
- 555 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
556 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
557 *ence on machine learning*, pp. 1861–1870. Pmlr, 2018.
- 558 Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of
559 the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary*
560 *computation*, 11(1):1–18, 2003.
- 561 Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane,
562 Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz,
563 et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous*
564 *Agents and Multi-Agent Systems*, 36(1):26, 2022.
- 565 Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and
566 Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping. *Advances*
567 *in Neural Information Processing Systems*, 33:15931–15941, 2020.
- 568 Rustam Issabekov and Peter Vamplew. An empirical comparison of two common multiobjective
569 reinforcement learning algorithms. In *AI 2012: Advances in Artificial Intelligence: 25th Aus-*
570 *tralian Joint Conference, Sydney, Australia, December 4-7, 2012. Proceedings 25*, pp. 626–636.
571 Springer, 2012.
- 572 Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for
573 exploration in reinforcement learning. In *International Conference on Machine Learning*, pp.
574 22594–22613. PMLR, 2023.
- 575 Haoye Lu, Daniel Herman, and Yaoliang Yu. Multi-objective reinforcement learning: Convexity,
576 stationarity and pareto optimality. In *The Eleventh International Conference on Learning Repre-*
577 *sentations*, 2022.
- 578 Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayara-
579 man, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via
580 coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- 581 Marlos C Machado, Marc G Bellemare, and Michael Bowling. Count-based exploration with the
582 successor representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol-
583 *ume 34*, pp. 5125–5133, 2020.
- 584 Gabriel B Margolis and Pulkit Agrawal. Walk these ways: Tuning robot control for generalization
585 with multiplicity of behavior. In *Conference on Robot Learning*, pp. 22–31. PMLR, 2023.
- 586 Farzan Memarian, Wonjoon Goo, Rudolf Lioutikov, Scott Niekum, and Ufuk Topcu. Self-supervised
587 online reward shaping in sparse-reward environments. In *2021 IEEE/RSJ International Confer-*
588 *ence on Intelligent Robots and Systems (IROS)*, pp. 2369–2375. IEEE, 2021.

- 594 Mayank Mittal, Calvin Yu, Qinxu Yu, Jingzhou Liu, Nikita Rudin, David Hoeller, Jia Lin Yuan,
595 Ritvik Singh, Yunrong Guo, Hammad Mazhar, Ajay Mandlekar, Buck Babich, Gavriel State,
596 Marco Hutter, and Animesh Garg. Orbit: A unified simulation framework for interactive robot
597 learning environments. *IEEE Robotics and Automation Letters*, 8(6):3740–3747, 2023. doi:
598 10.1109/LRA.2023.3270034.
- 599 Hossam Mossalam, Yannis M Assael, Diederik M Roijers, and Shimon Whiteson. Multi-objective
600 deep reinforcement learning. *arXiv preprint arXiv:1610.02707*, 2016.
- 601
602 Sriraam Natarajan and Prasad Tadepalli. Dynamic preferences in multi-criteria reinforcement learn-
603 ing. In *Proceedings of the 22nd international conference on Machine learning*, pp. 601–608,
604 2005.
- 605
606 Simone Parisi, Matteo Pirota, Nicola Smacchia, Luca Bascetta, and Marcello Restelli. Policy gra-
607 dient approaches for multi-objective sequential decision making. In *2014 International Joint*
608 *Conference on Neural Networks (IJCNN)*, pp. 2323–2330. IEEE, 2014.
- 609
610 Simone Parisi, Matteo Pirota, and Jan Peters. Manifold-based multi-objective policy search with
611 sample reuse. *Neurocomputing*, 263:3–14, 2017.
- 612
613 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration
614 by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787.
615 PMLR, 2017.
- 616
617 Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement.
618 In *International conference on machine learning*, pp. 5062–5071. PMLR, 2019.
- 619
620 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
621 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
622 *Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/
20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 623
624 Mathieu Reymond, Eugenio Bargiacchi, and Ann Nowé. Pareto conditioned networks. *arXiv*
625 *preprint arXiv:2204.05036*, 2022.
- 626
627 Mathieu Reymond, Conor F Hayes, Denis Steckelmacher, Diederik M Roijers, and Ann Nowé.
628 Actor-critic multi-objective reinforcement learning for non-linear utility functions. *Autonomous*
Agents and Multi-Agent Systems, 37(2):23, 2023.
- 629
630 Diederik M Roijers, Denis Steckelmacher, and Ann Nowé. Multi-objective reinforcement learn-
631 ing for the expected utility of the return. In *Proceedings of the Adaptive and Learning Agents*
workshop at FAIM, volume 2018, 2018.
- 632
633 Reuven Rubinstein. The cross-entropy method for combinatorial and continuous optimization.
634 *Methodology and computing in applied probability*, 1(2):127–190, 1999.
- 635
636 Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using
637 massively parallel deep reinforcement learning. In *Proceedings of the 5th Conference on Robot*
638 *Learning*, volume 164 of *Proceedings of Machine Learning Research*, pp. 91–100. PMLR, 2022.
639 URL [https://proceedings.mlr.press/v164/
rudin22a.html](https://proceedings.mlr.press/v164/rudin22a.html).
- 640
641 Manuela Ruiz-Montiel, Lawrence Mandow, and José-Luis Pérez-de-la Cruz. A temporal difference
642 method for multi-objective reinforcement learning. *Neurocomputing*, 263:15–25, 2017.
- 643
644 Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the
645 human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):
646 148–175, 2015.
- 647
648 Umer Siddique, Paul Weng, and Matthieu Zimmer. Learning fair policies in multi-objective (deep)
649 reinforcement learning with average and discounted rewards. In *International Conference on*
Machine Learning, pp. 8905–8915. PMLR, 2020.

- 648 Haoran Tang, Rein Houthoofd, Davis Foote, Adam Stooke, OpenAI Xi Chen, Yan Duan, John Schul-
649 man, Filip DeTurck, and Pieter Abbeel. # exploration: A study of count-based exploration for
650 deep reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
651
- 652 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,
653 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard
654 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- 655 Kristof Van Moffaert and Ann Nowé. Multi-objective reinforcement learning using sets of pareto
656 dominating policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.
657
- 658 Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforce-
659 ment learning: Novel design techniques. In *2013 IEEE symposium on adaptive dynamic pro-
660 gramming and reinforcement learning (ADPRL)*, pp. 191–199. IEEE, 2013a.
- 661 Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Hypervolume-based multi-objective
662 reinforcement learning. In *Evolutionary Multi-Criterion Optimization: 7th International Con-
663 ference, EMO 2013, Sheffield, UK, March 19-22, 2013. Proceedings 7*, pp. 352–366. Springer,
664 2013b.
- 665 Yuqing Xie, Chao Yu, Hongzhi Zang, Feng Gao, Wenhao Tang, Jingyi Huang, Jiayu Chen, Botian
666 Xu, Yi Wu, and Yu Wang. Multi-uav behavior-based formation with static and dynamic obstacles
667 avoidance via reinforcement learning. *arXiv preprint arXiv:2410.18495*, 2024.
668
- 669 Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective
670 reinforcement learning and policy adaptation. *Advances in neural information processing systems*,
671 32, 2019.
- 672 Chen Bo Calvin Zhang, Zhang-Wei Hong, Aldo Pacchiano, and Pulkit Agrawal. Orso: Ac-
673 celerating reward design via online reward selection and policy optimization. *arXiv preprint
674 arXiv:2410.13837*, 2024a.
- 675 Yihua Zhang, Prashant Khanduri, Ioannis Tsaknakis, Yuguang Yao, Mingyi Hong, and Sijia Liu.
676 An introduction to bilevel optimization: Foundations and applications in signal processing and
677 machine learning. *IEEE Signal Processing Magazine*, 41(1):38–59, 2024b.
678
- 679 Yunzhi Zhang, Pieter Abbeel, and Lerrel Pinto. Automatic curriculum learning through value dis-
680 agreement. *Advances in Neural Information Processing Systems*, 33:7648–7659, 2020.
681
- 682 Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On learning intrinsic rewards for policy gradient
683 methods. *Advances in neural information processing systems*, 31, 2018.
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A PRACTICAL GUIDELINES FOR DEPLOYING MORSE

A.1 GENERAL SUGGESTIONS

When deploying MORSE, we identify several important conditions which we recommend practitioners to follow.

- Make sure the inner-loop policy converges before performing outer-loop updates (including gradient updates and exploration). For bi-level gradient to hold, the policy must reach its local optimum, otherwise the gradient would be distorted. As for exploration, it relies on the policy’s current performance to decide whether to reset weights.
- Since we approximate the Hessian numerically, monitoring its stability is crucial. If the series diverges, the update should be skipped.
- Use slower outer-loop updates. Slowing down outer-loop updates, especially in the early stage of training, improves stability.
- Although MORSE enhances exploration of the reward weight space, extremely high-dimensional settings remain challenging. Stronger priors that narrow the weight range can significantly improve performance. If the search space is large with only a small valid region, more training iterations are required for sufficient exploration.
- Resetting the policy ensures high action entropy under new reward weights, enabling effective exploration.

A.2 APPLYING MORSE TO REAL WORLD ROBOTIC TRAINING

Domain randomization is widely used to improve sim-to-real transfer of RL policies, but the increased randomness it brings often makes training difficult. Here, we show how to integrate MORSE into a practical training workflow for real-world robotic control.

Rather than applying MORSE directly in a highly randomized environment, we adopt a two-stage training pipeline. Stage 1 performs reward shaping under minimal randomization, which allows MORSE to quickly search for appropriate reward weights. Stage 2 then trains the control policy in a fully domain-randomized environment using the reward function obtained from Stage 1.

We evaluate this pipeline on the IsaacLab quadcopter task. In Stage 1, we randomize only the target position and run MORSE 8 times to obtain 8 different weights. This stage repeats the Quadcopter experiment in Sec. 3. In Stage 2, we add mass and initial pose randomization to the environment, and train RL policies from scratch using both the MORSE-derived weights and the hand-crafted oracle reward. Each reward configuration is trained with PPO for 150 iterations over four seeds.

We present the results in Fig. 5b. Under domain randomization, the reward weights obtained by MORSE achieve performance comparable with the oracle reward. This demonstrates that MORSE-derived reward functions generalize effectively to more challenging training conditions and therefore can be applied to practical robotics settings.

B EXPERIMENTAL SETTINGS

B.1 TASK DEFINITION

In Sec. 6.1, we design three classes of synthetic functions to simulate optimization landscapes of varying difficulty: (1) SmoothPolynomial, a continuously differentiable random polynomial; (2) FixedNN, a function generated by a fixed randomly initialized neural network, which is non-smooth but exhibits no abrupt value shifts; (3) RandomSpiky, a smooth function with sparse sharp peaks. All functions are normalized to produce outputs within the range $[0,1]$. We present visualizations of representative instances from each class in Fig. 6.

We provide the task definition for RL environments in Tab. 3. In MuJoCo environments, the agent gets a task reward of 1 when it satisfies all task criteria, which are determined using the following protocol. First, we train a policy to convergence using the default reward weights provided in the original environment. Then, we perform multiple rollouts and record the minimum, maximum, and

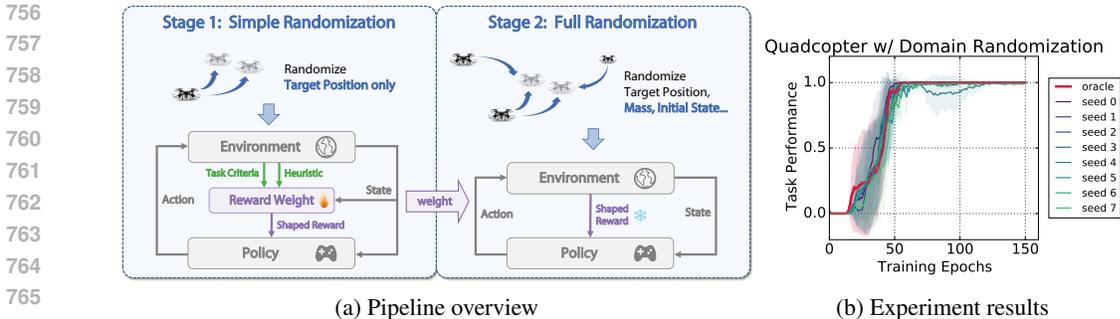


Figure 5: (a) 2-stage pipeline for MORSE under domain randomization. (b) Comparison between the oracle reward weight and MORSE-derived reward weights under domain randomization. MORSE reward achieves comparable performance to the oracle reward.

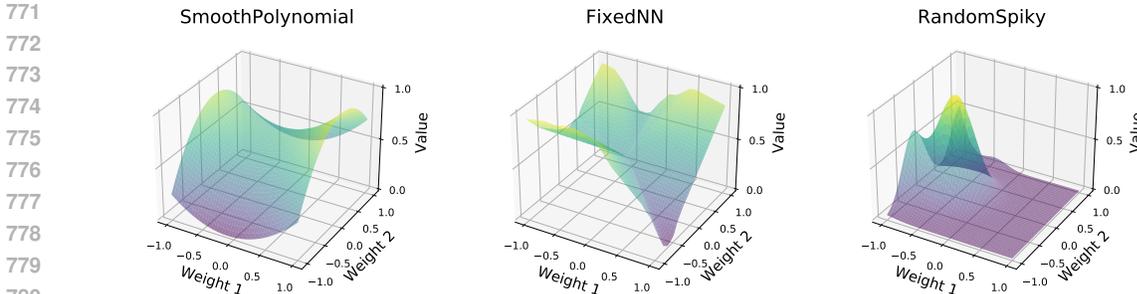


Figure 6: Example optimization landscapes.

mean values of each objective across episodes. Based on these statistics, we define task criteria such that the default policy achieves a success rate exceeding 90%.

To accelerate training, we shorten the episode length and modify the Unitree-A1 task from 9 objectives to 2 objectives. The maximum episode lengths for the environments are 100, 100, 100, 100, 200, 250, respectively.

B.2 HYPERPARAMETER

We use REINFORCE as the inner-loop optimizer for the CartPole experiments in Sec.4, and PPO for the MuJoCo environments and the Isaac Sim environments in Sec.6.2. Most hyperparameters follow the defaults of their respective implementations. Key hyperparameters are summarized in Table 4.

For CartPole, we use a single environment and set the maximum buffer size to 10,000. The inner-loop REINFORCE actor uses a learning rate of 0.001 and an L2 regularization coefficient of 0.25. The outer loop employs 5-step Neumann approximation with a learning rate of 0.001.

For MuJoCo, we use 32 parallel environments. Each PPO update uses a buffer size of 16,384 environment steps. The actor learning rate is 0.0003, value function coefficient is 0.5, maximum gradient is capped at 1.0, and the entropy coefficient is set to 0.

For Isaac Sim, we use 50 parallel environments. We use adaptive learning schedule for inner-loop policy training. The value function coefficient is 1.0, kl coefficient is 0.01, and the entropy coefficient is set to 0.01. The learning rate for Unitree-A1 is 1e-3, and the learning rate for Quadcopter is 5e-4.

The outer loop applies a softmax-based sampling strategy with temperature $\tau = 10$, where the sampling probability is given by $P(V_{\text{perf}_i}) = \frac{e^{V_{\text{perf}_i} * \tau}}{\sum_j e^{V_{\text{perf}_j} * \tau}}$. To preserve human-designed reward heuristics, we set the learning rate for the reward-weight network to 0.0005 and the coefficient learning

Environment	Robot	Task Criteria	Heuristic Functions	Weight
Cartpole		Final position > 0.5	Survival reward: 1 Position reward: $\mathbf{1}_{\text{pos} > 0.5}$ Interference reward: $\mathbf{1}_{\text{action}_x = \text{action}_{\text{mid}}}$	- - -
HalfCheetah		Final position > 5 control cost ≥ -0.6 for each step	Forward reward: $\frac{dx}{dt}$ Control cost: $- \text{action} _2^2$	0.4 1
Hopper		Final position > 0.8 control cost ≥ -0.004 for each step	Survival reward: 1 Forward reward: $\frac{dx}{dt}$ Control cost: $- \text{action} _2^2$	1 1 0.002
Walker2d		Final position > 0.8 control cost ≥ -0.006 for each step	Survival reward: 1 Forward reward: $\frac{dx}{dt}$ Control cost: $- \text{action} _2^2$	1 1 0.002
Quadcopter		L2 distance to goal < 4 linear velocity < 100 angular velocity < 10000	Position reward: $1 - \tanh(p - p^* ^2/0.8)$ Linear velocity cost: $-\Sigma v ^2$ Angular velocity cost: $-\Sigma \omega ^2$	15 0.05 0.01
Unitree-A1-2obj		L2 tracking error < 0.22 L2 action rate < 40	Tracking reward: $e^{-(v-v^*)^2}$ Action rate penalty: $- a_t - a_{t-1} _2^2$	1.5 0.01
Unitree-A1-9obj		L2 tracking error < 0.22 L2 action rate < 40	Linear xy velocity error: $e^{-(v_{xy}-v_{xy}^*)^2}$ Angular yaw velocity error: $e^{-(v_{yaw}-v_{yaw}^*)^2}$ Linear z velocity penalty: $- v_z _2^2$ Angular-xy velocity penalty: $- \omega_{xy} _2^2$ Joint torque penalty: $- \tau _2^2$ Joint acceleration penalty: $- \ddot{q} _2^2$ Action rate penalty: $- a_t - a_{t-1} _2^2$ Air-time reward: $\Sigma_i (A_i - \tau) F_i, \mathbf{H} u_{cmd} > 0.1$ Non-flat orientation penalty: $- g_{b,xy} _2^2$	1.5 0.75 -2 -0.05 -0.0002 -2.5e-07 -0.01 0.25 -2.5
Reach-Franka		Position tracking error < 0.2 Fine-grained tracking reward > 0.5 Orientation tracking error < 0.6 L2 action rate < 100	Position tracking error: $ p - p^* _2$ Fine-grained pos reward: $1 - \tanh(p - p^* /0.1)$ Orientation tracking error: $d_{\text{quat}}(\mathbf{q}, \mathbf{q}^*)$ Action rate penalty: $- a_t - a_{t-1} _2^2$	-0.2 0.1 -0.1 -0.0001
Lift-Franka		L2 object position tracking error < 0.22 L2 action rate < 25000	Reaching object reward: $1 - \tanh\left(\frac{ p_{\text{obj}} - p_{\text{obj}}^* _2}{0.1}\right)$ Lifting object reward: $\text{lift} = \mathbf{1}_{\{z_{\text{obj}} > 0.04\}}$ Object goal tracking: $\text{lift} * \left[1 - \tanh\left(\frac{ p^* - p_{\text{obj}} _2}{0.3}\right)\right]$ Fine-grained goal tracking: $\text{lift} * \left[1 - \tanh\left(\frac{ p^* - p_{\text{obj}} _2}{0.05}\right)\right]$ Action rate penalty: $- a_t - a_{t-1} _2^2$	1.0 15.0 16.0 5.0 -0.0001

Table 3: Task definition for RL environments.

Environment	Inner Loop n_epochs	Outer Loop	
		gradient frequency	reset frequency
Cartpole	15	1	-
HalfCheetah	10	5	25
Hopper	10	5	25
Walker2d	20	15	50
Quadcopter	5	20	100
Unitree-A1-2obj	20	20	250
Unitree-A1-9obj	20	20	250
Reach-Franka	20	20	250
Lift-Franka	20	50	500

Table 4: Hyperparameters for RL environments.

rate to 0.0025. The number of gradient updates for reward weights is adaptive, with a minimum of 3 and a maximum of 10 iterations per outer-loop step, terminating early upon convergence. Similarly, the prediction network is trained using a dynamic number of iterations, from 1 to 1000, to allow it to overfit known data points during each update cycle.

Each run takes less than 1.5 hour on a 4090 GPU.

C ADDITIONAL EXPERIMENTS AND ANALYSIS

C.1 ALTERNATIVE DESIGN CHOICES OF MORSE

MORSE is a modular framework, and each component can be flexibly instantiated.

For inner-loop policy training, users can choose any RL algorithms they like, such as SAC(Haarnoja et al., 2018).

For the approximation of outer-loop gradient, while we adopt the implicit function approach following (Gupta et al., 2023), users can apply other meta-learning approaches, such as (Hu et al., 2020). We compare both approaches in MuJoCo-hard environments under the Gradient Only setting as in Sec. 3. As shown in Fig. 7, the two estimators yield similar performance, likely because both are trapped in local optima in the absence of exploration.

MORSE also permits diverse novelty metrics for exploration. Sec. 4 reports results with random sampling, and Fig. 8 further extends the comparison to include Bayesian Optimization (Shahriari et al., 2015). Across these variants, RND achieves the strongest performance, providing more informative exploration signals for reward-weight search.

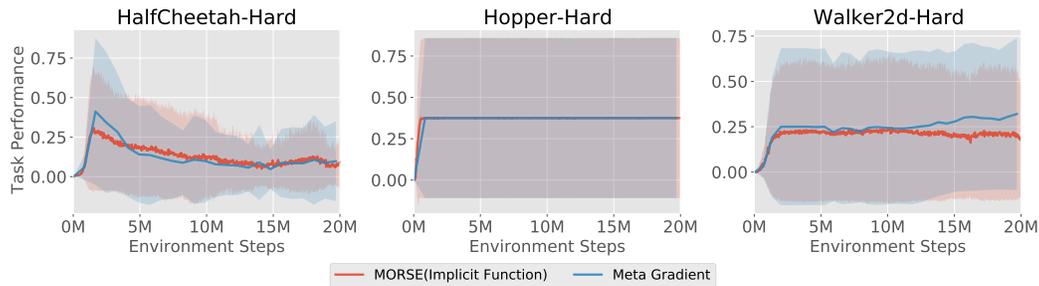


Figure 7: Comparison of implicit-function and meta-gradient estimators for approximating the outer-loop gradient.

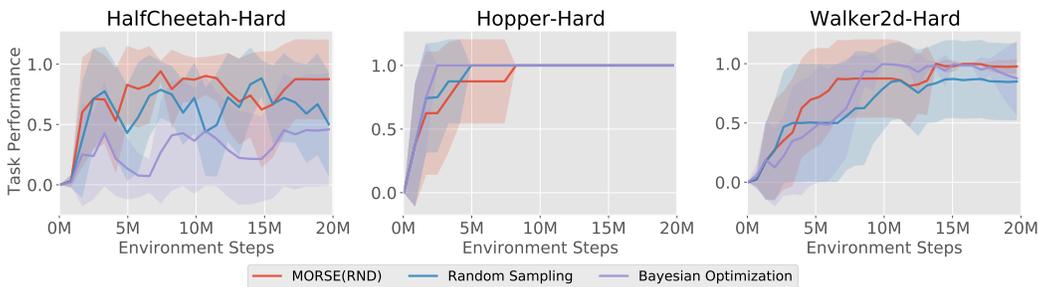


Figure 8: Comparison of RND, random sampling, and Bayesian Optimization as exploration metrics.

C.2 OPTIMIZATION PATH FOR PRELIMINARY VALIDATION

Following the experiment setting in Sec. 6.1.1, we visualize the optimization path of each experiment in Fig. 9. The results show that MORSE can quickly find the most novel reward weight and perform a stable gradient update near the reward weight that performs well, resulting in optimal performance.

C.3 SENSITIVITY ANALYSIS OF OUTER-LOOP PARAMETERS

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

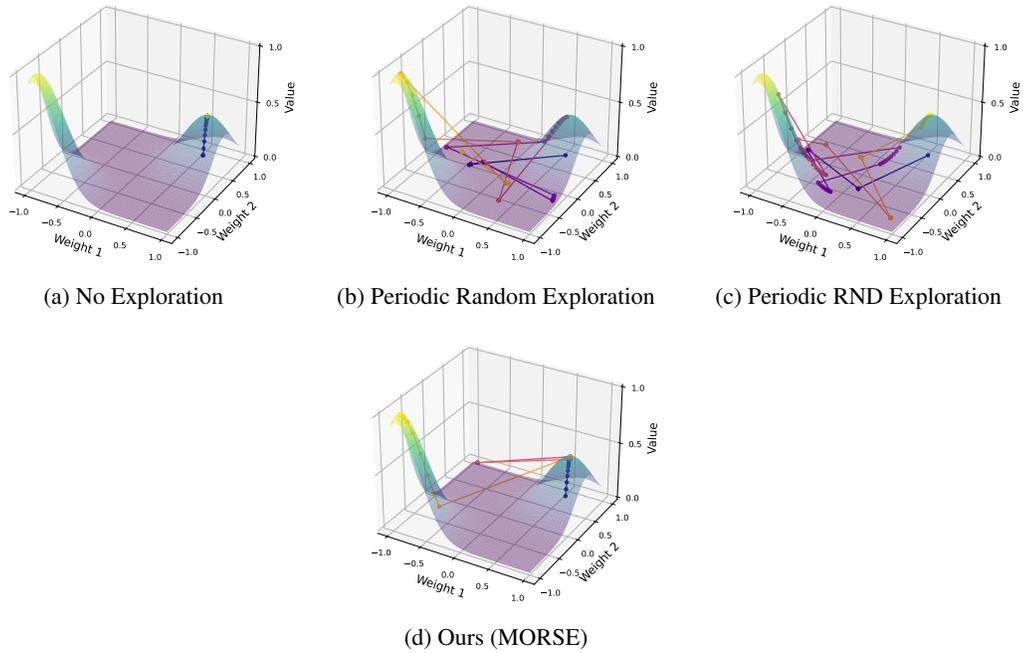


Figure 9: Visualization of optimization path. MORSE optimizes towards the global optimum on the lower left.

The outer loop of MORSE has two principal hyperparameters, the bi-level gradient learning rate and the outer-loop exploration interval. Using the same experimental setup as in Sec. 3, we sweep each hyperparameter independently and report the results in Fig. 10 and Fig. 11.

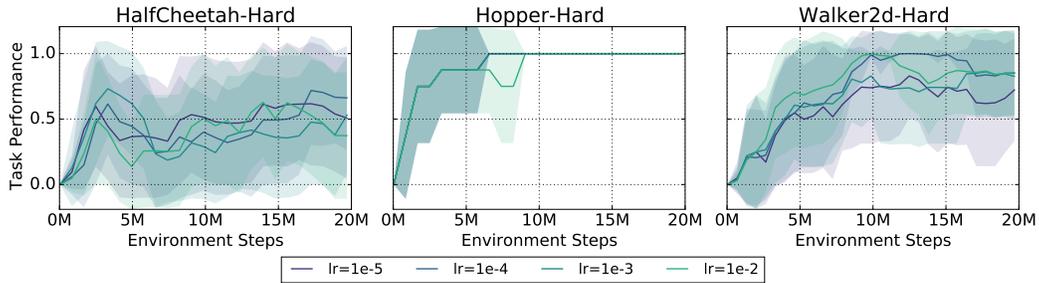


Figure 10: Sensitivity to outer-loop learning rate.

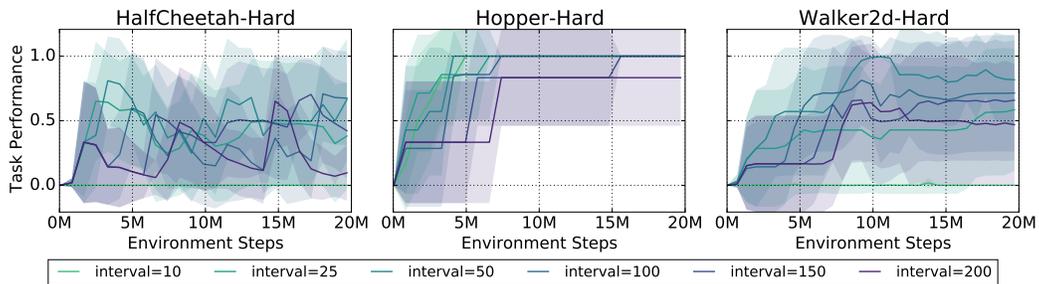


Figure 11: Sensitivity to exploration interval $T_{explore}$.

Our key observations are as follows. First, the overall policy performance is relatively insensitive to the change of learning rate across the tested range, probably due to our constraints on the update process. We clip the reward within its valid range, and use gradient clipping to prevent aggressive updates.

Second, exploration interval can greatly influence performance. When the interval is too small, the outer loop triggers exploration before the inner-loop policy adapts to the current reward weights, causing promising weights to be incorrectly judged as poor solutions and thus discarded. However, when the interval becomes too large, the outer loop can only explore a limited number of times within a fixed training budget, which reduces the chance of discovering high-quality weights. Therefore, we recommend carefully choosing the explore frequency and extending the training budget, so that the inner loop policy has enough time to converge to the current reward before the next exploration, and the outer loop can perform sufficient exploration to find the optimal weight.

C.4 PCA VISUALIZATION OF WEIGHT SELECTION

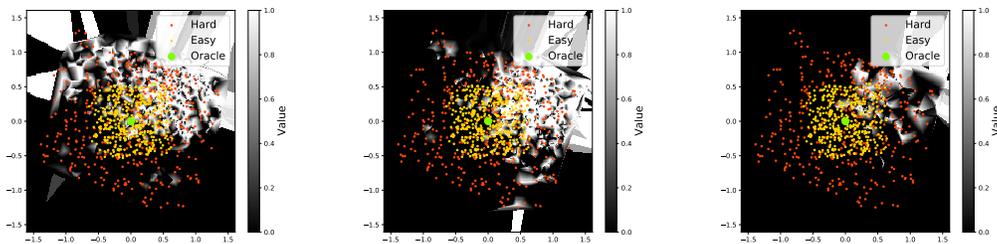


Figure 12: Visualization of the optimization planes of MuJoCo environments using PCA. The color denotes task success rate. The brighter, the higher. From left to right: HalfCheetah, Hopper, Walker2d.

In this section, we visualize the changes of the reward weight with each outer-loop update during the training process in MuJoCo environments. Since the reward weight is three-dimensional, we apply PCA to reduce it to two dimensions for visualization.

In Fig. 12, we illustrate the relationship between reward weight and task success rate across the entire reward weight space. We collect reward weights and success rates from all experiments and then apply interpolation and edge filling to complete the missing regions. We also highlight the human-tuned oracle reward weight and mark the valid reward weight ranges for both easy and hard settings with 400 randomly sampled points. The resulting optimization landscape exhibits multiple local optima, consistent with our hypothesis that robotic task optimization landscapes are often rugged and therefore challenging for gradient-based optimization.

Fig. 13 presents the evolution of reward weight for MuJoCo-Hard environments. We compare four methods: **MORSE**, **Gradient**, and two ablation variants, namely, **wo/ Performance Metric** and **wo/ Novelty Metric**. All methods start optimizing from the same initial point.

Gradient is highly sensitive to the choice of initial point. When the initial point falls to undesired areas, gradient updates can be ineffective. **wo/ Performance Metric** escapes the optimal solution. **wo/ Novelty Metric** covers the reward weight space more slowly. For more challenging tasks such as Walker2d, which have lower exploration frequency and fewer reset chances, **wo/ Novelty Metric** limits the ability to find suitable reward weights.

D THE USE OF LARGE LANGUAGE MODELS (LLMs)

We use LLMs for code debugging and paper polishing.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

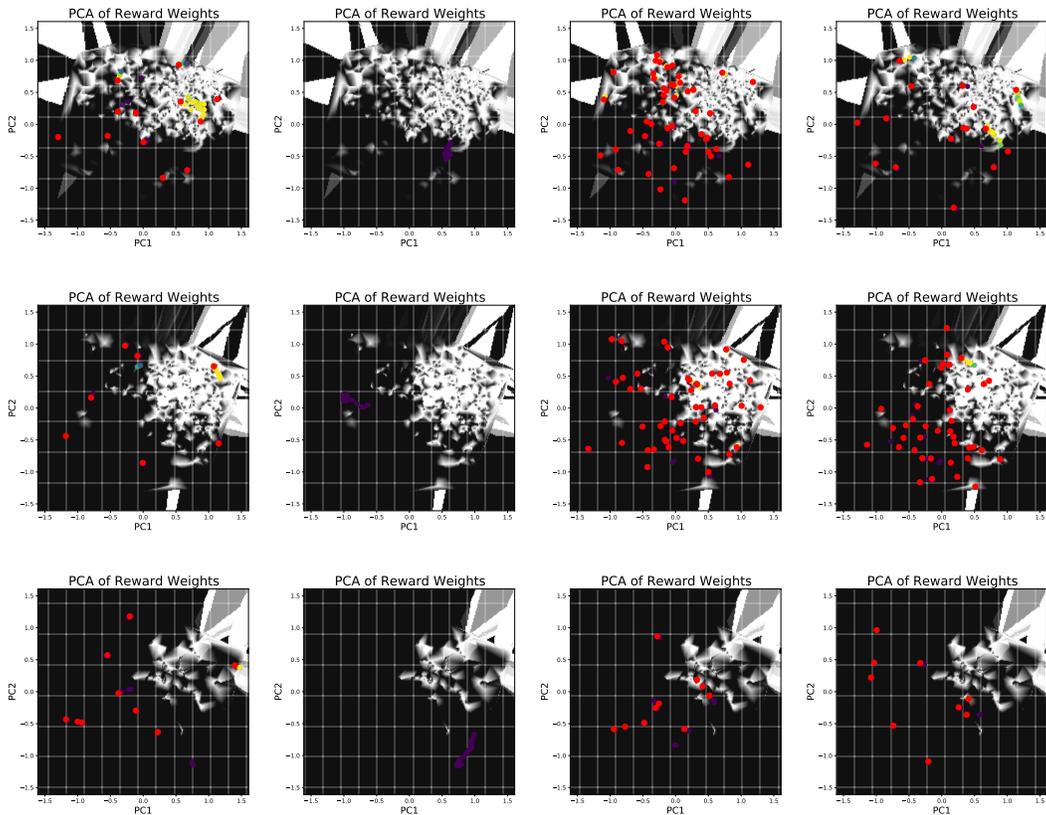


Figure 13: Evolution of reward weights in MuJoCo-Hard environments. Each row corresponds to a different environment, and each column represents a method. From top to bottom: HalfCheetah-Hard, Hopper-Hard, Walker2d-Hard. From left to right: **MORSE**, **Gradient**, **wo/ Performance Metric**, **wo/ Novelty Metric**. The color of the scatter points indicates the number of outer-loop updates, with lighter colors representing more updates. Red points denote new starting points selected when resetting the reward weight.