FROM REWARD SHAPING TO Q-SHAPING: ACHIEVING UNBIASED LEARNING WITH LLM-GUIDED KNOWL-EDGE

Anonymous authors

Paper under double-blind review

Abstract

Q-shaping is an extension of Q-value initialization and serves as an alternative to reward shaping for incorporating domain knowledge to accelerate agent training, thereby improving sample efficiency by directly shaping Q-values. This approach is both general and robust across diverse tasks, allowing for immediate impact assessment while guaranteeing optimality. We evaluated Q-shaping across 20 different environments using a large language model (LLM) as the heuristic provider. The results demonstrate that Q-shaping significantly enhances sample efficiency, achieving an 16.87% average improvement across the 20 tasks compared to the best baseline, and a 226.67% improvement compared to LLM-based reward shaping methods. These findings establish Q-shaping as an effective and unbiased alternative to conventional reward shaping in reinforcement learning.

1 INTRODUCTION

Reinforcement learning (RL) can solve complex tasks but often faces sample inefficiency. For example, AlphaGo (Silver et al., 2016) required approximately 4 weeks of training on 50 GPUs, learning from 30 million expert Go game positions to reach a 57% accuracy. Similarly, training a real bipedal soccer robot required 9.0×10^8 environment steps, amounting to 68 hours of wall-clock time for the full 1v1 agent (Haarnoja et al., 2024). These cases demonstrate the significant computational demands of RL.



Figure 1: Agent behavior across different algorithms. "Vanilla" refers to traditional RL algorithms,
 "reward shaping" refers to reward shaping-enhanced RL algorithms, and "Q-shaping" refers to
 Q-shaping-enhanced RL algorithms. Q-shaping impacts agent behavior quickly, enabling rapid
 evolution and improvement in the quality of heuristic functions. In contrast, reward shaping requires
 extensive training time before the impact of the heuristic reward becomes apparent.

To improve efficiency, popular methods include (1) imitation learning, (2) residual reinforcement
learning, (3) reward shaping, and (4) Q-value initialization. Yet, each has limitations: imitation
learning requires expert data (Garg et al., 2021; Chang et al., 2024; Kostrikov et al., 2020), residual
RL needs a well-designed controller (Johannink et al., 2019; Trumpp et al., 2023), and Q-value
initialization (Nakamoto et al., 2024) demands precise estimates. Therefore, reward shaping (Xie

et al.; Ma et al., 2023) is the most practical approach, as it avoids the need for expert trajectories or
 predefined controllers.

Reward shaping methods fall into two main categories: (1) potential-based reward shaping 057 (PBRS) (Ng et al., 1999) and (2) non-potential-based reward shaping (NPBRS) (Ng et al., 1999). 058 PBRS provides state-based heuristic rewards and ensures that optimality is preserved by following the potential function rule, as defined by . NPBRS, on the other hand, refers to reward shaping 060 methods that do not adhere to the potential function rule, and as a result, the learned policy does not 061 guarantee optimality. Additionally, reward shaping methods often suffer from a slow verification 062 process, requiring completion of training to assess the impact of the heuristic reward, which limits 063 their development, as noted by Ma et al. (2023). Lastly, designing high-quality reward functions 064 remains a challenging and often frustrating task for researchers, hindering the adoption of these methods (Ma et al., 2023). 065

With the growing popularity of large language models (LLMs), LLM-guided reinforcement learning (RL) has emerged as a promising field. This approach leverages the strong understanding capabilities of LLMs to guide RL agents in exploration or policy updates. Existing research has focused on two main areas: LLM-based policy generation and LLM-guided reward design. For example, Chen et al. (2021); Micheli et al. (2022) utilize LLMs to enhance policy decisions, while Kwon et al. (2023); Carta et al. (2023); Ma et al. (2023) employ LLMs to design reward structures. Although these works have improved task success rates, the challenges associated with reward shaping remain unresolved.

In this work, we introduce a novel framework, Q-shaping, which leverages domain knowledge from large language models (LLMs) to guide agent exploration. Q-shaping offers two key advantages over reward shaping:

077 078

079

081

1. **Remain Optimality**: Q-shaping inspires exploration by modifying Q-values during training while ensuring that the agent's optimality remains unaffected upon convergence.

- 2. Efficient Heuristic Verification: Unlike reward shaping methods, which require waiting until the end of training to observe the impact of the reward heuristic, Q-shaping enables experimenters to verify and refine heuristic guidance rapidly during training.
- Figure 1 illustrates the agent behavior across different algorithms.

In the "Q-shaping Framework" section, we present theoretical analysis and proofs demonstrating that
Q-shaping preserves optimality while using imprecise Q-values to improve exploration and sample
efficiency. In the experimental section, we use GPT-40 as a heuristic provider and compare Q-shaping
with popular baselines, achieving an average improvement of 16.87% across 20 tasks. Compared to
LLM-guided reward shaping methods like T2R (Xie et al.) and Eureka (Ma et al., 2023), Q-shaping
achieves up to 226.67% improvement in episodic total rewards while enhancing task success rates.

089 090

2 RELATED WORK

091 092 093

094 095 2.1 HEURISTIC REINFORCEMENT LEARNING

There are four common approaches to incorporating domain knowledge into reinforcement learning to enhance sample efficiency: (1) Imitation Learning, (2) Residual Policy, (3) Reward Shaping, and (4) Q-value Initialization.

Imitation Learning requires access to expert trajectories, as demonstrated by works such as GAIL (Ho
& Ermon, 2016), where agents learn by mimicking expert behavior. However, the reliance on highquality expert data limits its applicability in complex tasks. Residual Policy (Johannink et al., 2019)
methods involve designing a controller to guide agent actions, but this manual design process restricts
their scalability and generality.

Q-value initialization, although promising, often requires precise Q-value estimates to derive an
effective policy. For instance, Cal-QL (Nakamoto et al., 2024) employs calibrated Q-values to
enhance agent exploration, but these calibrated values still rely on expert knowledge, making Q-value
design more challenging than reward shaping. Consequently, few studies have pursued this direction
due to the inherent difficulty in obtaining accurate Q-values compared to reward shaping.

108 Reward shaping directly modifies the reward function to influence agent behavior, improving training 109 efficiency without requiring expert trajectories or manual controller design. This approach has 110 been refined to address diverse learning scenarios, such as in Inverse Reinforcement Learning (IRL) 111 (Ziebart et al., 2008; Wulfmeier et al., 2015; Finn et al., 2016) and Preference-based RL (Christiano 112 et al., 2017; Ibarz et al., 2018; Lee et al., 2021; Park et al., 2022). Additionally, various heuristic techniques have been introduced, including unsupervised auxiliary task rewards (Jaderberg et al., 113 2016), count-based reward heuristics (Bellemare et al., 2016; Ostrovski et al., 2017), and self-114 supervised prediction error heuristics (Pathak et al., 2017; Stadie et al., 2015; Oudeyer & Kaplan, 115 2007). 116

However, reward shaping often suffers from inaccuracies in the heuristic functions and a slow verification process, which limits its effectiveness in certain applications.

119 120

121

2.2 LLM\VLM AGENT

122 LLMs/VLMs can achieve few-shot or even zero-shot learning in various contexts, as demonstrated 123 by works such as Voyager (Wang et al., 2023), ReAct (Yao et al., 2022), SLINVIT (Zhang et al., 124 2024), and SwiftSage (Lin et al., 2024). In the field of robotics, VIMA Jiang et al. (2022) employs 125 multimodal learning to enhance agents' comprehension capabilities. Additionally, the use of LLMs 126 for high-level control is becoming a trend in control tasks (Shi et al., 2024; Liu et al., 2023; Ouyang 127 et al., 2024). In web search, interactive agents (Gur et al., 2023; Shaw et al., 2024; Zhou et al., 2023) can be constructed using LLMs/VLMs. Moreover, frameworks have been developed to reduce 128 the impact of hallucinations, such as decision reconsideration (Yao et al., 2024; Long, 2023), self-129 correction (Shinn et al., 2023; Kim et al., 2024), and observation summarization (Sridhar et al., 130 2023). 131

132 133

134

2.3 LLM-ENHANCED RL

135 Relying on the understanding and generation capabilities of large models, LLM-enhanced RL has 136 become a popular field (Du et al., 2023; Carta et al., 2023). Researchers have investigated the diverse 137 roles of large models within reinforcement learning (RL) architectures, including their application in reward design (Kwon et al., 2023; Wu et al., 2024; Carta et al., 2023; Chu et al., 2023; Yu et al., 138 2023; Ma et al., 2023), information processing (Paischer et al., 2022; 2024; Radford et al., 2021), as a 139 policy generator, and as a generator within large language models (LLMs) (Chen et al., 2021; Micheli 140 et al., 2022; Robine et al., 2023; Chen et al., 2022). While LLM-assisted reward design has improved 141 task success rates (Ma et al., 2023; Xie et al.), it often introduces bias into the original Markov 142 Decision Process (MDP) or fails to provide sufficient guidance for complex tasks. Additionally, the 143 verification process is time-consuming, which slows down the pace of iterative improvements. 144

145 146

147

3 NOTATION

For policy definition, the space of all possible policies is denoted as Π. A policy π : S → Δ(A)
defines a conditional distribution over actions given states. A deterministic policy μ : S → A is
a special case of π, where one action is selected per state with a probability of 1. We define an
"activity matrix" A^π ∈ ℝ^{S×Z} for each policy, encoding π's state-conditional state-action distribution.
Specifically, A^π(s, ⟨s, a⟩) := π(a|s) if s = s, otherwise A^π(s, ⟨s, a⟩) := 0. The value function is
defined as v: Π → S → ℝ or q: Π → S × A → ℝ, both with bounded outputs. The terms q and v
represent discrete matrix representations, where v(s) and q(s, a) specifically denote the outputs of

An optimal policy for an MDP \mathcal{M} , denoted by $\pi_{\mathcal{M}}^*$, is one that maximizes the expected return under the initial state distribution: $\pi_{\mathcal{M}}^* := \arg \max_{\pi} \mathbb{E}_{\rho}[\mathbf{v}_{\mathcal{M}}^{\pi}]$. The state-wise expected returns of this optimal policy are represented by $\mathbf{v}_{\mathcal{M}}^{\pi_{\mathcal{M}}^*}$. The Bellman consistency equation for the MDP \mathcal{M} at \mathbf{x} is given by $\mathcal{B}_{\mathcal{M}}(\mathbf{x}) := \mathbf{r} + \gamma P \mathbf{x}$. Notably, $(\mathbf{v}_{\mathcal{M}}^{\pi})^*$ is the unique vector that satisfies $(\mathbf{v}_{\mathcal{M}}^{\pi})^* = A^{\pi} \mathcal{B}_{\mathcal{M}}((\mathbf{v}_{\mathcal{M}}^{\pi})^*)$. We abbreviate \mathbf{q}^* as $(\mathbf{q}_{\mathcal{M}}^{\pi_{\mathcal{M}}^*})^*$ for some MDP ξ .

Datasets We define fundamental concepts essential for fixed-dataset policy optimization. Let $D := \{\langle s, a, r, s' \rangle\}^d$ represent a dataset of d transitions. From this dataset, we can construct a local MDP \mathcal{D} and derive a local optimal Q-value function, denoted as q_D^* .

172 Within the Q-shaping framework, let $\hat{\mathbf{q}}$ denote the Q-function learned from TD estimation and 173 Q-shaping. The LLM outputs are categorized into two types: goodQ, which encourages exploration, 174 and badQ, which discourages it. Let $G_{LLM} := \{(s, a, Q) \mid Q > 0\}^d$ represent the dataset of d175 heuristic pairs focused on encouraging agent exploration. Similarly, $B_{LLM} := \{(s, a, Q) \mid Q \le 0\}^d$ 176 denotes the dataset of d heuristic pairs aimed at preventing exploration. The complete collection of 177 LLM outputs is given by $D_{LLM} := \{G_{LLM}, B_{LLM}\}.$

Convergence An agent is considered to have converged when it reaches 80% of the peak performance. The peak performance is defined as the highest performance achieved by any of the baseline methods.

4 Q-SHAPING FRAMEWORK

178

179

181 182

183

190

191

192 193

194 195 196

200

208

In the Q-learning framework, an experience buffer D is used to store transitions from the Markov Decision Process (MDP), supporting both online and offline training. To estimate the Q-values for (s, a) pairs, the Temporal-Difference (TD) update method leverages this experience buffer. The Q-function derived from the trained Q-values determines the policy by maximizing $q(s, \cdot)$, making accurate Q-value estimation crucial for policy quality and effective exploration.

To enhance exploration, Q-shaping integrates both the experience buffer and heuristic guidance from a large language model (LLM) into the Q-value estimation process. The **Heuristic TD Update**, which defines this Q-shaping process, is given by:

 $\hat{\mathbf{q}}^{k+1}(s,a) \leftarrow \begin{cases} \hat{\mathbf{q}}^k(s,a) + \alpha h(s,a), & \text{if } (s,a) \in D_{LLM}^k \setminus \mathcal{D}, \\ \hat{\mathbf{q}}^k(s,a) + \alpha (\hat{\mathbf{q}}_{TD}^k(s,a) + h(s,a)), & \text{if } (s,a) \in D_{LLM}^k \cap \mathcal{D}. \end{cases}$

where $\hat{\mathbf{q}}_{TD}^k(s, a)$ represents the temporal-difference (TD) update estimation of $\mathbf{q}(s, a)$ at step k, expressed as: $\hat{\mathbf{q}}_{TD}^k(s, a) = r(s, a, s') + \gamma \hat{\mathbf{q}}^k(s, a)$. Here, D_{LLM}^k denotes the set of (s, a, h(s, a))pairs provided by the LLM at iteration k.

With this formulation, the Heuristic Bellman Optimal Operator can be expressed as:

$$\hat{\mathbf{q}}^{k+1}(s,a) = \mathcal{T}_h \hat{\mathbf{q}}^k(s,a) \tag{1}$$

$$= r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} \hat{\mathbf{q}}^k(s',a') + h(s,a), \quad (s,a) \in D^k_{LLM} \cap \mathcal{D}.$$
(2)

4.1 UNBIASED OPTIMALITY

The Q-value represents a high-level abstraction of an agent's interaction with the environment. It encapsulates the expected cumulative reward by integrating critical elements such as rewards r, transition probabilities P, states s, actions a, and the policy π . Changes in any of these components directly affect the Q-values.

SAC (Haarnoja et al., 2018) and MCTS (Browne et al., 2012) use action-bonus heuristics to enhance
 training efficiency but risk biasing the learned policy away from optimality. In contrast, Q-shaping,
 supported by Theorem 1, enhances learning with heuristic guidance while ensuring convergence to
 the optimal Q-values of the local MDP.

Theorem 1 (Contraction and Convergence of $\hat{\mathbf{q}}$). Let \mathcal{T}_h be the heuristic Bellman operator for the sampled MDP \mathcal{D} , and let $\gamma \in [0,1)$ be the discount factor. The operator \mathcal{T}_h satisfies the following contraction property in the metric space $(\mathcal{X}, \|\cdot\|_{\infty})$:

222

223

224

225

226

227 228

229 230

231 232

233

239

248

 $\|\mathcal{T}_h(\hat{\mathbf{q}}) - \mathcal{T}_h(\hat{\mathbf{q}}')\|_{\infty} \leq \gamma \|\hat{\mathbf{q}} - \hat{\mathbf{q}}'\|_{\infty},$

where $\hat{\mathbf{q}}, \hat{\mathbf{q}}' \in \mathcal{X}$ are any two value functions. Thus, \mathcal{T}_h is a γ -contraction operator.

As a result, repeated applications of the heuristic Bellman operator through the heuristic Temporal Difference (TD) update,

 $\hat{\mathbf{q}} \leftarrow \mathcal{T}_h(\hat{\mathbf{q}}),$

will converge to the unique fixed point $\hat{\mathbf{q}}_{\mathcal{D}}^*$. Furthermore, since $\hat{\mathbf{q}}$ and \mathbf{q} are updated on the same MDP and Follow Assumption A.2, the following equivalence holds:

 $\hat{\mathbf{q}}_{\mathcal{D}}^* = \mathbf{q}_{\mathcal{D}}^*.$

Proof. See Appendix A.2

4.2 UTILIZING IMPRECISE Q VALUE ESTIMATION

At the early training stage, the Q-values for different actions are nearly identical, leading the policy to execute actions randomly. To address this, we leverage the LLM's domain knowledge to provide positive Q-values for actions that contribute to task success and negative Q-values for actions that do not. The imprecise Q-values provided by the LLM can be categorized into two types: overestimations and underestimations.

Underestimation of Non-Optimal Actions An agent does not need to fully traverse the entire
 state-action space to identify the optimal trajectory that leads to task success. Therefore, imprecise
 Q-value estimation can be effectively utilized to guide the agent's exploration.

For instance, consider a scenario where the agent is required to control a robot arm to operate on a drawer located in front of it. In this case, actions such as moving the arm backward or upward are evidently unhelpful in finding the optimal trajectory. Assigning very low Q-values to these non-contributory actions discourages the agent from exploring them, thereby enhancing sample efficiency.

249 **Overestimation of Near-Optimal Ac-**250 tions At the initial training phase (iteration step k = 0, let action a 251 be assumed to have the highest es-252 timated Q-value for a given state s, 253 while a^* denotes the true optimal action. 254 This assumption leads to the inequal-255 ity $\hat{\mathbf{q}}(s, a^*) < \hat{\mathbf{q}}(s, a) < \mathbf{q}^*(s, a^*)$. 256 Consequently, the agent is predisposed 257 to explore actions around the subop-258 timal a in its search for states, given 259 that $\mu(s) = \max_a \hat{\mathbf{q}}(s, \cdot) + \epsilon$, where 260 $\epsilon \sim \mathcal{N}(0, \delta^2)$.

261 However, the number of steps required 262 to discover the optimal action a^* is in-263 herently constrained by the environment 264 and the distance between a and a^* . To 265 expedite this exploration process, we in-266 troduce an action a_{LLM} suggested by Algorithm 1 Q-shaping

- 1: **Require**: Good Q-set G_{llm} , Bad Q-set B_{llm} provided by the LLM, RL solver A
- 2: Goal: Compute the average performance over 10 runs
- 3: Initialize: Start 20 agents { $Agent_1, Agent_2, \dots, Agent_{20}$ }
- 4: # for each agent, do:
- 5: agent.explore(steps = 5000)
- 6: # Apply Q-shaping and Policy-shaping
- 7: agent.q_shaping(G_{llm}, B_{llm})
- 8: agent.policy_shaping(G_{llm}, B_{llm})
- 9: # Further exploration
- 10: agent.explore(steps = 10000)
- 11: # Synchronize agents
- 12: agent.wait()
- 13: # Remove 10 lower-performing agents
 14: agent.remove if latter()
- 15: # Continued exploration and training
- 16: agent.explore_and_train()
- 17: **Output**: Average performance over 10 runs

the LLM, replacing a via Q-shaping guided by the loss function in Equation 3 to enhance sample efficiency. Given the assumption $|a_{LLM} - a^*| < |a - a^*| < \delta$, we can express $\mu(s) = a_{LLM} + \epsilon$. Consequently, the agent has a higher chance of selecting a^* , significantly improving the likelihood of identifying the optimal trajectory. In conclusion, by letting the LLM provide the *goodQ* set and *badQ* set, the agent is guided to prioritize exploring actions suggested by the LLM, thereby enhancing sample efficiency. Over time, as indicated by Hasselt (2010); Fujimoto et al. (2018) and Theorem 1, \hat{q} converges towards the locally optimal Q-function. We now present the theoretical upper bound on the sample complexity required for $\hat{\mathbf{q}}$ to converge to $\mathbf{q}_{\mathcal{D}}^*$ for any given MDP \mathcal{D} :

Theorem 2 (Convergence Sample Complexity). The sample complexity n required for $\hat{\mathbf{q}}$ to converge to the local optimal fixed-point \mathbf{q}_D^* with probability $1 - \delta$ is:

$$n > \mathcal{O}\left(\frac{|S|^2}{2\epsilon^2} \ln \frac{2|S \times A|}{\delta}\right)$$

Proof. See proof at A.4.

Theorem 2 establishes an upper bound on the sample complexity, indicating that the imprecise Q-values provided by the LLM will be corrected within a finite number of steps. Therefore, any heuristic values can be introduced during the early training iterations, and the Q-shaping framework will adapt to inaccurate Q-values over time.

4.3 Algorithm Implementation

For the implementation of Q-shaping, we employ TD3 (Fujimoto et al., 2018) as the RL solver (backbone) and GPT-40 as the heuristic provider, introducing three additional training phases: (1) Q-Network Shaping (2) Policy-Network Shaping, and (3) High-performance agent selection. Pseudo-code 1 outlines the detailed steps of the Q-shaping framework.

O-Network Shaping In the Q-shaping framework, the LLM is tasked with providing a set of (s, a, Q) pairs to guide exploration. This approach is particularly crucial during the early training stage when it is challenging for the agent to independently discover expert trajectories. Traditional RL solvers often require a substantial number of steps to identify the correct path to success, leading to sample inefficiency. The goal of the Q-shaping framework is to leverage the provided (s, a, Q)pairs to accelerate exploration and help the agent quickly identify the optimal path.



Figure 2: Q-shaping prompt. There is a general code template that specifies the required structure for the generated code. In addition to the template, three key pieces of information are necessary to generate an effective heuristic function: the code template, an introduction to the environment provided in the paper, and the environment configuration file.

To obtain D_{LLM} , we construct a general code template as the prompt as illustrated in Figure 2, supplemented by task-specific environment configuration files and a detailed definition of the observation and action spaces within the simulator. Subsequently, we apply the loss function $L_{q-shaping}$ to update the Q-function:

$$L_{q-shaping}(\theta) = E_{(s_i, a_i, Q_i) \sim D_g} (Q_i - \hat{\mathbf{q}}_{\theta}(s_i, a_i))^2$$
(3)

Policy-Network Shaping In most reinforcement learning (RL) algorithms, the policy is derived from the Q-function, where the policy is optimized to execute actions that maximize the Q-value given a state. The policy update is expressed as: $\mu(s) = \arg \max_a \hat{\mathbf{q}}(s, \cdot)$

While introducing a learning rate and target policy can help stabilize the training process and prevent fluctuations in the policy network, this approach often slows down the convergence speed. To accelerate this adaptation, we introduce a "Policy-Network Shaping" phase designed to allow the policy to quickly align with the good actions and avoid the bad actions provided by the LLM.

332 The shaping loss function is defined as:

333 334

341 342

343

344

345

346

347 348

349 350

351 352

353

354

355

364 365

366

367

368

369

370

371

372

373

374

375

$$L_{policy-shaping} = \lambda_1 \mathbb{E}_{(s,a)\sim G_{LLM}} \left[\|\mu(s) - a\|^2 \right] - \lambda_2 \mathbb{E}_{(s,a)\sim B_{LLM}} \left[\|\mu(s) - a\|^2 \right]$$
(4)

335 , where $(s, a) \sim G_{LLM}$ and $(s, a) \sim B_{LLM}$ represent state-action pairs sampled from the LLM-336 provided *goodQ* and *badQ* sets, respectively, and λ_1 and λ_2 are hyperparameters controlling the 337 influence of the LLM-guided shaping.

With this "Policy-Network Shaping" phase, researchers can quickly observe the impact of heuristic values, facilitating the rapid evolution of heuristic quality, ultimately leading to a more efficient exploration process and faster convergence to optimal behavior.

High-Performance Agent Selection With Q-network shaping and policy-network shaping, the policy is initialized to perform actions suggested by the LLM. However, due to the randomness of initial states and the learning process, the shaped agent may still perform poorly and have a chance of failing to learn effectively. To address this, each agent is allowed 10,000 steps to test its performance after the Q-network and policy network initialization. Following this evaluation, weaker agents are removed, and only the top-performing agent continues with the training process.

5 EXPERIMENT SETTINGS

We investigate the following **hypotheses** through a series of four experiments:

- 1. Q-shaping can enhance sample efficiency in reinforcement learning.
- 2. Q-shaping can adapt to incorrect or hallucinated heuristics while maintaining optimality.
- 3. Q-shaping outperforms LLM-based reward shaping methods.
- 4. LLM can design heuristic functions that provide s, a, Q altogether.



Figure 3: Evaluation environments span a diverse set of robot types and tasks, ranging from simple pendulum systems to humanoid control. The 20 tasks cover a variety of state dimensions, robotic types, and reward structures

To validate these hypotheses, we conducted three primary experiments and two ablation study. GPT-40 served as the heuristic provider, while TD3 was employed as the reinforcement learning (RL) backbone, forming **LLM-TD3**. As illustrated in Figure 3, Q-shaping and various baseline methods were evaluated across 20 distinct tasks involving drones, robotic arms, and other robotic control challenges. Below, we describe the specific experiments and their objectives:

- 1. **Sample Efficiency Experiment:** This experiment compares Q-shaping with four baseline methods to evaluate its impact on sample efficiency during training.
- Comparison with LLM-based Reward Shaping: Q-shaping, which integrates domain knowledge to assist in agent training, is compared with Text2Reward and Eureka to evaluate its performance relative to existing LLM-based reward shaping approaches.

- 3. LLM Quality Evaluation: Although Q-shaping guarantees optimality, its reliance on LLM-provided heuristics may influence performance. This experiment evaluates the quality of different LLM outputs.
- 4. **Ablation Study on Q-shaping phases:** Q-shaping introduces three key training phases. This experiment isolates and examines the contribution of each phase to overall performance.
- 5. **Teachability Experiment:** This experiment evaluates the teachability of different LLMs by analyzing how few interactions can improve code quality and performance.



Figure 4: Learning curve comparison of each algorithm across 20 tasks.

Environments We evaluate Q-shaping and baselines across 20 distinct environments, including 8 from Gymnasium Classic Control and MuJoCo (Todorov et al., 2012), 9 from MetaWorld (Yu et al., 2020), and 3 from PyFlyt (Tai et al., 2023). Notably, the robotic arm and drone environments used are less commonly studied, making it unlikely that the LLM was pretrained on these specific environments.

Baselines For the sample efficiency experiments, we compared Q-shaping against several baseline algorithms, including CleanRL-PPO, CleanRL-SAC (Huang et al., 2022), DDPG (Lillicrap et al., 2015), and TD3 (Fujimoto et al., 2018). When evaluating Q-shaping against other reward shaping methods, we selected Text2Reward and Eureka as baselines. In the LLM-type ablation study, we assessed the performance of different LLMs: o1-Preview, GPT-4o-Mini, Gemini-1.5-Flash (Team et al., 2023), DeepSeek-V2 (DeepSeek-AI et al., 2024), and Yi-Large (Young et al., 2024). - **Text2Reward**: Text2Reward leverages GPT-4 to generate reward functions from natural language task descriptions. In this study, we use provided prompts to describe the MetaWorld tasks, with SAC as the baseline RL algorithm for training policies.

411 - **Eureka**: Eureka utilizes an evolutionary algorithm to iteratively evolve reward functions based on 412 task performance, refining the reward function over successive generations to improve task success 413 rates. In this work, K (iteration batch size) is set to 8, and N = 5 (search iterations) is used. We use 414 GPT-40 as the reward generator and CleanRL-PPO as the backbone reinforcement learning algorithm. 415 The prompts used to generate reward functions are detailed in Appendix B.3.





432 **Metrics** To evaluate sample efficiency, we measure the number of steps required to reach 80% of 433 peak performance, where peak performance is defined as the highest performance achieved by any 434 baseline agent. For clarity in visualization, improvements exceeding 150% are truncated to 150%. 435

Each algorithm is tested 10 times, and the average evaluation performance is reported. Evaluations 436 are conducted at intervals of 5,000 steps. During each evaluation, the agent is tested over 10 episodes, 437 and the average episodic return is plotted to form the learning curve. 438

In this study, we do not specify a fixed seed for each run. Using a fixed seed results in a unique initial state when the environment is reset, which simplifies learning and makes it challenging to accurately verify the effectiveness and generalization capabilities of each algorithm.

441 442 443

444 445

446

447

448

449 450

451

452

453

454

455

456

457

458 459

460

461

439

440

6 **RESULTS AND ANALYSIS**

O-Shaping Outperforms Best Baseline by an Average of 16.87% Across 20 Tasks As shown in Figure 5 and Figure 4, Q-shaping demonstrated a notable improvement over both the best baseline and TD3 across 20 tasks. On average, Q-shaping improves performance by 16.87% compared to the best baseline and by 55.39% compared to TD3, highlighting its effectiveness in enhancing sample efficiency and task performance. This supports H1.



Figure 6: Learning curve comparison between Q-shaping and LLM-based reward shaping methods. The evaluation was conducted on four Meta-World environments: door-close, drawer-open, window*close*, and *sweep-into*, with peak performance serving as the basis for comparison.

462 Q-Shaping Outperforms LLM-Based Reward Shaping Methods by 226.67% Q-shaping 463 achieved substantial improvements over both the Eureka and T2R baselines, as shown in Figure 6. 464 The comparison is based on peak performance across the evaluated Meta-World environments. 465

Compared to the best baseline, LLM-466 TD3 improved by 5.16% in the 467 door-close task, 81.89% in drawer-468 open, 715.67% in window-close, and 469 103.96% in sweep-into, resulting in an 470 average peak performance improve-471 ment of 226.67%. 472

Eureka LLM-TD3

LLM-based reward shaping methods, 473

though capable of improving task suc-474

cess rates (Ma et al., 2023; Xie et al.), often bias optimality and, as shown in Table 1, require 475 substantial time to evaluate the effectiveness of reward heuristics. In contrast, Q-shaping achieves a 476 226.67% improvement over the best LLM-based reward shaping methods and requires only a few 477 steps to validate the heuristic function. This supports H2 and H3.

Algorithm

- 478 Most LLMs Can Provide Correct
- 479 Heuristic Functions We evaluated 480 the quality of LLM-generated heuris-481 tic functions from five perspectives: 482 (1) adherence to the required code 483 template, (2) correctness of the assigned Q-values, (3) accuracy of the 484 state-action dimension, (4) complete-485 ness of the generated code, and (5)

Table 2: Evaluation of LLM Quality in Outputting Heuristic Values

Table 1: Additional training steps required to derive

across four Meta-World environments.

 8×10^6

 $1.5 imes 10^3$

door-close-v2 drawer-open-v2

effective heuristic functions for LLM-TD3 and Eureka

 8×10^{6}

 2×10^3

sweep-into-v2

 8×10^{6}

 3×10^{3}

window-close-v2

 8×10^6

 2×10^3

Metric	o1-Preview	GPT-40	Gemini	DeepSeek-V2.5	yi-large
Template Adherence (%)	100.0	100.0	40.0	100.0	100.0
Correct Q-values (%)	100.0	100.0	60.0	100.0	100.0
Correct State-Action Dim (%)	100.0	100.0	80.0	100.0	100.0
Code Completeness (%)	100.0	100.0	20.0	100.0	100.0
Bug-Free (%)	100.0	100.0	20.0	100.0	100.0
Average (%)	100.0	100.0	44.0	100.0	100.0

presence of bugs in the generated code. Each LLM was prompted 10 times with the same request, and we quantified their performance using a correctness rate across these metrics.

Table 3: Ablation Study on Additional Training Phases. The study evaluates the impact of three key training phases—Q-Network Shaping, Policy-Network Shaping, and Agent Selection—across four Meta-World environments: *door-close*, *drawer-open*, *window-close*, and *sweep-into*. Effectiveness is measured by convergence steps (\downarrow) , with "Failed" indicating algorithms that did not reach the convergence threshold within 10^6 steps.

	Phase			Environ		
Q-shaping	Policy-shaping	Selection	door-close-v2	drawer-open-v2	sweep-into-v2	window-close-v2
×	×	×	Failed	Failed	Failed	759999
\checkmark	×	×	Failed	310000	Failed	570000
×	\checkmark	×	30000	340000	Failed	215000
\checkmark	\checkmark	×	30000	275000	860000	195000
\checkmark	\checkmark	\checkmark	25000	265000	790000	165000

502Correctness of the assigned Q-values means that state-
action pairs (s, a) from the LLM-generated goodQ set504must be assigned Q-values greater than zero, while those505from the badQ set must be assigned Q-values less than or506equal to zero.

The results, as shown in Table 2, indicate that most LLMs, including o1-Preview, GPT-40, DeepSeek-V2.5, and yilarge, provided correct heuristic functions with a 100% success rate across all evaluation metrics. However, Gemini exhibited poorer performance, achieving only 44% on average. This supports H4.





Figure 7: Teachability of different LLMs. The x-axis represents the number of interactions, while the y-axis shows the average performance across four tasks: *Door-Close, Drawer-Open, Sweep-Into*, and *Window-Close*.

Few Interactions Significantly Improve Code Quality Figure 7 illustrates the teachability of LLMs within the Q-shaping framework. Remarkably, all models achieved high performance within just 3 to 4 interactions, suggesting that the primary issue with the initial generated code lies in parameter tuning rather than structural flaws.

7 CONCLUSION

We propose Q-shaping, an alternative framework that integrates domain knowledge to enhance sample efficiency in reinforcement learning. In contrast to traditional reward shaping, Q-shaping offers two key advantages: (1) it preserves optimality, and (2) it allows for rapid verification of the agent's behavior. These features enable experimenters or LLMs to iteratively refine the quality of heuristic values without concern for the potential negative impact of poorly designed heuristics. Experimental results demonstrate that Q-shaping significantly improves sample efficiency and outperforms LLM-guided reward shaping methods across various tasks.

We hope this work encourages further research into advanced techniques that leverage LLM outputs to guide and enhance the search process in reinforcement learning.

540 REFERENCES

- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
 Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves
 Oudeyer. Grounding large language models in interactive environments with online reinforcement
 learning. In *International Conference on Machine Learning*, pp. 3676–3713. PMLR, 2023.
- Jonathan D Chang, Dhruv Sreenivas, Yingbing Huang, Kianté Brantley, and Wen Sun. Adversarial
 imitation learning via boosting. *arXiv preprint arXiv:2404.08513*, 2024.
- 555
 556
 557
 Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 modeling. Advances in neural information processing systems, 34:15084–15097, 2021.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Kun Chu, Xufeng Zhao, Cornelius Weber, Mengdi Li, and Stefan Wermter. Accelerating reinforce ment learning of robotic manipulations via feedback from large language models. *arXiv preprint arXiv:2311.02379*, 2023.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024. URL https://arxiv.org/abs/2405.04434.
- Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek
 Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language
 models. In *International Conference on Machine Learning*, pp. 8657–8677. PMLR, 2023.
- 575
 576
 576
 577
 Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pp. 49–58. PMLR, 2016.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn: Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34: 4028–4039, 2021.
- Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and
 Aleksandra Faust. A real-world webagent with planning, long context understanding, and program
 synthesis. *arXiv preprint arXiv:2307.12856*, 2023.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- Tuomas Haarnoja, Ben Moran, Guy Lever, Sandy H Huang, Dhruva Tirumala, Jan Humplik, Markus
 Wulfmeier, Saran Tunyasuvunakool, Noah Y Siegel, Roland Hafner, et al. Learning agile soccer
 skills for a bipedal robot with deep reinforcement learning. *Science Robotics*, 9(89):eadi8022, 2024.

594 595 596	Hado Hasselt. Double q-learning. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 23. Curran Associates, Inc., 2010. URL https://proceedings.neurips.cc/paper_files/paper/
597	2010/file/091d584fced301b442654dd8c23b3fc9-Paper.pdf.
598	Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. Advances in neural
599	information processing systems, 29, 2016.
601	Shangui Huang, Dausslan Formand Julian Dassa, Chang Va, Jaff Praga, Dinam Chakrabartu, Kinal
602	Mehta and Ioão G M Araúio Cleantl' High-quality single-file implementations of deep rein-
603	forcement learning algorithms. <i>Journal of Machine Learning Research</i> , 23(274):1–18, 2022. URL
604	http://jmlr.org/papers/v23/21-1342.html.
605	Boria Ibarz, Jan Leike, Tobias Poblen, Geoffrey Irving, Shane Legg, and Dario Amodei, Reward
606	learning from human preferences and demonstrations in atari. Advances in neural information
607	processing systems, 31, 2018.
608	Man Indenkana Maladaman Maik Waisingh Manian Carmonaki Tam Sakard Ind 71 sika David
609	Silver and Koray Kayukcuoglu Reinforcement learning with unsupervised auxiliary tasks. In
610	International Conference on Learning Representations, 2016.
612	
613	Yuntan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei,
614	prompts In NeurIPS 2022 Foundation Models for Decision Making Workshop 2022
615	prompts. In rearing 5 2022 Foundation inducts for Decision inducing workshop, 2022.
616	Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll,
617	Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for
618	IEFE 2019
619	
620	Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks.
621	Advances in Neural Information Processing Systems, 36, 2024.
622	Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution
623	matching. In International Conference on Learning Representations, 2020.
625	Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh, Reward design with language
626	models. arXiv preprint arXiv:2303.00001, 2023.
627	Visin Log Lours Smith and Distor Abbael, Dabbles Feedback officient interactive minforcement
628	learning via relabeling experience and unsupervised pre-training arXiv preprint arXiv:2106.05091
629	2021.
630	
631	Inmothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Dean Wierstra. Continuous control with dean reinforcement lograming. arXiv
632	preprint arXiv:1509.02971. 2015.
633	
034	Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula,
636	Prinviraj Ammanabroiu, Yejin Unoi, and Alang Ken. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. Advances in Neural Information Processing Systems
637	36, 2024.
638	
639	Huihan Liu, Alice Chen, Yuke Zhu, Adith Swaminathan, Andrey Kolobov, and Ching-An Cheng.
640	Interactive robot learning from verbal correction. <i>arxiv preprint arxiv:2510.17555</i> , 2025.
641	Jieyi Long. Large language model guided tree-of-thought. arXiv preprint arXiv:2305.08291, 2023.
642	Yecheng Jason Ma William Liang Guanzhi Wang De-An Huang Oshert Rastani Dinesh Javaraman
643	Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding
644	large language models. In The Twelfth International Conference on Learning Representations,
645	2023.
046 647	Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. <i>arXiv preprint arXiv:</i> 2209.00588, 2022.

648 Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral 649 Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. 650 Advances in Neural Information Processing Systems, 36, 2024. 651 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: 652 Theory and application to reward shaping. In Icml, volume 99, pp. 278-287, 1999. 653 654 Georg Ostrovski, Marc G Bellemare, Aäron Oord, and Rémi Munos. Count-based exploration with 655 neural density models. In International conference on machine learning, pp. 2721–2730. PMLR, 656 2017. 657 Pierre-Yves Oudeyer and Frederic Kaplan. What is intrinsic motivation? a typology of computational 658 approaches. Frontiers in neurorobotics, 1:108, 2007. 659 660 Yutao Ouyang, Jinhan Li, Yunfei Li, Zhongyu Li, Chao Yu, Koushil Sreenath, and Yi Wu. Long-661 horizon locomotion and manipulation on a quadrupedal robot with large language models. arXiv 662 preprint arXiv:2404.05291, 2024. 663 Fabian Paischer, Thomas Adler, Vihang Patil, Angela Bitto-Nemling, Markus Holzleitner, Sebastian 664 Lehner, Hamid Eghbal-Zadeh, and Sepp Hochreiter. History compression via language models 665 in reinforcement learning. In International Conference on Machine Learning, pp. 17156–17185. 666 PMLR, 2022. 667 668 Fabian Paischer, Thomas Adler, Markus Hofmarcher, and Sepp Hochreiter. Semantic helm: A human-readable memory for reinforcement learning. Advances in Neural Information Processing 669 Systems, 36, 2024. 670 671 Jongjin Park, Younggyo Seo, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Surf: 672 Semi-supervised reward learning with data augmentation for feedback-efficient preference-based 673 reinforcement learning. In 10th International Conference on Learning Representations, ICLR 674 2022. International Conference on Learning Representations, 2022. 675 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration 676 by self-supervised prediction. In International conference on machine learning, pp. 2778–2787. 677 PMLR, 2017. 678 679 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 680 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 681 models from natural language supervision. In *International conference on machine learning*, pp. 8748-8763. PMLR, 2021. 682 683 Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models 684 are happy with 100k interactions. arXiv preprint arXiv:2303.07109, 2023. 685 686 Peter Shaw, Mandar Joshi, James Cohan, Jonathan Berant, Panupong Pasupat, Hexiang Hu, Urvashi 687 Khandelwal, Kenton Lee, and Kristina N Toutanova. From pixels to ui actions: Learning to follow instructions via graphical user interfaces. Advances in Neural Information Processing Systems, 36, 688 2024. 689 690 Lucy Xiaoyang Shi, Zheyuan Hu, Tony Z Zhao, Archit Sharma, Karl Pertsch, Jianlan Luo, Sergey 691 Levine, and Chelsea Finn. Yell at your robot: Improving on-the-fly from language corrections. 692 arXiv preprint arXiv:2403.12910, 2024. 693 Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic 694 memory and self-reflection. arXiv preprint arXiv:2303.11366, 2023. 695 696 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, 697 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016. 699 Abishek Sridhar, Robert Lo, Frank F Xu, Hao Zhu, and Shuyan Zhou. Hierarchical prompting 700 assists large language model on web navigation. In The 2023 Conference on Empirical Methods in 701 Natural Language Processing, 2023.

- Bradly C Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- Jun Jet Tai, Jim Wong, Mauro Innocente, Nadjim Horri, James Brusey, and Swee King Phang. Pyflyt–
 uav simulation environments for reinforcement learning research. *arXiv preprint arXiv:2304.01305*, 2023.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. IEEE, 2012.
- Raphael Trumpp, Denis Hoornaert, and Marco Caccamo. Residual policy learning for vehicle control of autonomous racing cars. In *2023 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1–6. IEEE, 2023.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and
 Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv* preprint arXiv:2305.16291, 2023.
- Yue Wu, Yewen Fan, Paul Pu Liang, Amos Azaria, Yuanzhi Li, and Tom M Mitchell. Read and reap the rewards: Learning to play atari with the help of instruction manuals. *Advances in Neural Information Processing Systems*, 36, 2024.
- Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforce ment learning. *arXiv preprint arXiv:1507.04888*, 2015.
- Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang, and Tao Yu. Text2reward: Reward shaping with language models for reinforcement learning. In *The Twelfth International Conference on Learning Representations*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montserrat Gonzalez
 Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language
 to rewards for robotic skill synthesis. In *7th Annual Conference on Robot Learning*, 2023.
- 747
 748
 749
 749
 750
 740
 740
 741
 742
 743
 744
 744
 744
 745
 746
 747
 747
 748
 749
 749
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
 740
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
 Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building
 autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.
- Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.

A APPENDIX

758 A.1 ADDITIONAL NOTATION 759

Datasets. In practise, a batch of data will be sampled from a distribution $\Phi : Dist(\mathcal{Z})$, which is the collected local MDP. A batch D' containing d tuples $\langle s, a, r, s' \rangle$ is sampled as $D' \sim \Phi_d$, where pairs $\langle s, a \rangle$ are drawn from Φ , and rewards r and subsequent states s' are sampled independently from the reward function $\mathcal{R}_D(\cdot|\langle s, a \rangle)$ and the transition function $P_D(\cdot|\langle s, a \rangle)$, respectively.

764 Given a dataset or batch D, we denote $D(\langle s, a \rangle)$ as the multi-set of all $\langle r, s' \rangle$ pairs, and use $\ddot{\mathbf{n}}_D \in$ 765 $\mathbb{R}^{|\mathcal{Z}|}$ to denote the count vector, where $\ddot{\mathbf{n}}_D(\langle s, a \rangle) := |D(s, a)|$. We define the empirical reward 766 vector as $\mathbf{r}_D(\langle s, a \rangle) := \sum_{r,s' \in D(\langle s, a \rangle)} \frac{r}{|D(\langle s, a \rangle)|}$ and empirical transition matrix as $P_D(s'|\langle s, a \rangle) :=$ 767 $\sum_{r,s'\in D(\langle s,a\rangle)} \frac{\mathbb{I}(s'=s')}{|D(\langle s,a\rangle)|} \text{ for all state-action pairs with } \ddot{\mathbf{n}}_D(\langle s,a\rangle) > 0. \text{ For state-action pairs where}$ 768 $\ddot{\mathbf{n}}_D(\langle s, a \rangle) = 0$, the maximum-likelihood estimates of reward and transition cannot be clearly defined, 769 so they remain unspecified. The bounds hold no matter how these values are chosen, so long as \mathbf{r}_D is 770 bounded and P_D is stochastic. The empirical policy of a dataset D is defined as $\hat{\pi}_D(a|s) := \frac{|D(\langle s, a \rangle)|}{|D(\langle s, \cdot \rangle)|}$ 771 772 except where $\ddot{\mathbf{n}}_D(\langle s, a \rangle) = 0$, where it can similarly be any valid action distribution. The empirical 773 visitation distribution of a dataset D is computed analogously to the regular visitation distribution but uses P_D in place of P. Thus it's given by $\frac{1}{1-\gamma} (I - \gamma A^{\pi} P_D)^{-1}$. 774

Lemma 1 (Decomposition). For any MDP ξ and policy π , consider the Bellman fixed-point equation given by, let $(\mathbf{v}_{\xi}^{\pi})^*$ be defined as the unique value vector such that $(\mathbf{v}_{\xi}^{\pi})^* = A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}(\mathbf{v}_{\xi}^{\pi})^*)$, and let \mathbf{v} be any other value vector. Assume that $\pi(a|s) = 1$ if $a = \arg \max_{a}(\mathbf{q}_{\xi}^{\pi})^*(s, a)$, otherwise $\pi(a|s) = 0$. We have:

780 781

782

789

791

792 793 794

798 799 800

801 802

804 805

809

$$|\mathbf{q}_{\xi}^{*}(s,\mu(s)) - \mathbf{q}(s,\mu(s))| = |((I - \gamma A^{\pi}P_{\xi})^{-1}(A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}\mathbf{v}) - \mathbf{v}))(s)|$$
(5)

Proof.

$$A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}\mathbf{v}) - \mathbf{v} = A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}\mathbf{v}) - (\mathbf{v}_{\xi}^{\pi})^{*} + (\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v}$$

$$= A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}\mathbf{v}) - A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}(\mathbf{v}_{\xi}^{\pi})^{*}) + (\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v}$$

$$= \gamma A^{\pi} P_{\xi}(\mathbf{v} - (\mathbf{v}_{\xi}^{\pi})^{*}) + ((\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v})$$

$$= ((\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v}) - \gamma A^{\pi} P_{\xi}((\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v})$$

$$= (I - \gamma A^{\pi} P_{\xi})((\mathbf{v}_{\xi}^{\pi})^{*} - \mathbf{v})$$

Note that $(\mathbf{v}_{\xi}^{\pi})^* = A^{\pi}(\mathbf{q}_{\xi}^{\pi})^*$, After we expand the value function we have:

$$(I - \gamma A^{\pi} P_{\xi})^{-1} (A^{\pi} (\mathbf{r}_{\xi} + \gamma P_{\xi} \mathbf{v})) = A^{\pi} (\mathbf{q}_{\xi}^{\pi})^{*} - \mathbf{v}$$
$$= A^{\pi} (\mathbf{q}_{\xi}^{\pi})^{*} - A\mathbf{q}$$

By indexing at $\langle s, \mu(s) \rangle$, we have:

$$|\mathbf{q}_{\xi}^{*}(s,\mu(s)) - \mathbf{q}(s,\mu(s))| = |((I - \gamma A^{\pi}P_{\xi})^{-1}(A^{\pi}(\mathbf{r}_{\xi} + \gamma P_{\xi}\mathbf{v}) - \mathbf{v}))(s)|$$

Lemma 2 (Convergence Bound). Since that s' and r are sampled independently and identically distributed (iid) from $P_D(\cdot|s, a)$ and $R_D(\cdot|s, a)$ respectively. Let D' denotes the batch of data sample from D. Then, with probability at least $1 - \delta$, we have:

$$|\mathbf{q}_{D'}^*(s,\mu(s)) - \hat{\mathbf{q}}(s,\mu(s))| \le \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right)\sum_{s'}\nu_{D'}(s'|s_0=s)\frac{1}{\sqrt{\ddot{\mathbf{n}}_{D'}(\langle s',\mu(s')\rangle)}}$$

Proof. See proof at A.3

A.2 PROOF OF THEOREM 1

Assumption A.1. The heuristic h(s, a) provided by the LLM does not change with the training steps, i.e., $h^k(s, a) = h(s, a)$ for all k = 0, 1, 2, ...

Assumption A.2. The heuristic h(s, a) is only used during the initial training steps and is removed after some step k_0 , i.e., for all training steps $k \ge k_0$, the heuristic term is not provided.

Note that **q** is the matrix representation of the Q function. In the proof of this section, we use a more general $Q : \mathbb{R}^{\mathbb{Z}} \to \mathbb{R}$ to represent the Q function. The heuristic TD update for \hat{Q} iteration is:

$$\hat{Q}^{k+1}(s,a) = (1-\alpha)\hat{Q}^k(s,a) + \alpha \left(r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} \hat{Q}(s',a') + h(s,a) \right)$$

We can define a **Bellman optimal operator** \mathcal{T}_h based on the heuristic TD update as follows:

$$\hat{Q}^{k+1}(s,a) = \mathcal{T}_h \hat{Q}^k = r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} \hat{Q}^k(s',a') + h(s,a)$$

Suppose training framework Q-shaping satisfies assumption A.1. Then we prove that the Bellman optimal operator \mathcal{T}_h is γ -contraction operator on \hat{Q} :

$$\begin{aligned} \|\mathcal{T}_h \hat{Q} - \mathcal{T}_h \hat{Q}'\|_{\infty} &= \gamma \max_{s,a \in \mathcal{S}, \mathcal{A}} |\sum_{s'} P(s'|s,a) [\max_{a'} \hat{Q}(s',a') - \max_{a'} \hat{Q}'(s',a')]| \\ &\leq \gamma \max_{s,a \in \mathcal{S}, \mathcal{A}} |\max_{s'} |\left(\max_{a'} \hat{Q}(s',a') - \max_{a'} \hat{Q}'(s',a') \right)|| \\ &= \gamma \|\hat{Q} - \hat{Q}'\|_{\infty} \end{aligned}$$

The optimal Q-function for the new update formula, without assumption A.2, is defined as:

$$\hat{Q}^*(s,a) = r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} \hat{Q}^*(s',a') + h(s,a)$$

 \mathcal{T}_h is a γ -contraction operator on \hat{Q} . This means that as the number of iterations k increases, \hat{Q} will approach the heuristic fixed point, which is biased. Under assumption A.2, the heuristic TD update will degenerate into the TD update. Without the influence of the heuristic term, the Q-values will be estimated solely from the local MDP \mathcal{D} .

Next, we prove that the converged heuristic-guided Q function is equivalent to the traditional Q function. Define the following:

 Θ_H denotes the set of terminal states,

 Θ_{H-1} denotes the set of states one step before the terminal,

 Θ_1 denotes the set of states at the initial step.

For all $s \in \Theta_{H-1}$ and some action a, it is clear that $\hat{Q}^*(s, a) = Q^*(s, a)$, because:

$$Q^*(s,a)|_{s\in\Theta_{H-1}} = \hat{Q}^*(s,a)|_{s\in\Theta_{H-1}} = r(s,a) + \gamma \sum_{s'\in\Theta_H} \mathbf{1}_{s\in\Theta_H} \max_{a'} Q^*(s',a') = r(s,a)$$

For all $s \in \Theta_{H-2}$ and some action a, we have:

$$\hat{Q}^*(s,a)|_{s\in\Theta_{H-2}} = r(s,a) + \gamma \sum_{s'\in\Theta_{H-1}} P(s'|s,a) \max_{a'} \hat{Q}^*(s',a')$$

$$= r(s,a) + \gamma \sum_{s' \in \Theta_{H-1}} P(s'|s,a) \max_{a'} Q^*(s',a')$$

$$= r(s,a) + \gamma \sum_{s' \in \Theta_{H-1}} P(s'|s,a) \max_{a'} Q^*(s',a')$$

 $= Q^*(s,a)|_{s \in \Theta_{H-2}}$

With sufficient iterations, we have: $\hat{Q}^* = Q^*$. Specifically, we have: $\mathbf{q}^* = \hat{\mathbf{q}}^*$ for some MDP \mathcal{D} .

A.3 PROOF OF LEMMA 2

Let D' be a batch of data, and D denotes the replay buffer, consider that for any $\langle s, a \rangle$, the expression $\mathbf{r}_{D'}(\langle s, a \rangle) + \gamma P_{D'}(\langle s, a \rangle) \mathbf{v}^{\pi}$ can be equivalently expressed as an expectation of random variables,

$$\mathbf{r}_{D'}(\langle s, a \rangle) + \gamma P_{D'}(\langle s, a \rangle) \mathbf{v} = \frac{1}{\ddot{\mathbf{n}}_{D'}(\langle s, a \rangle)} \sum_{r, s' \in D'(\langle s, a \rangle)} r + \gamma \mathbf{v}(s')$$

each with expected value:

$$\mathbb{E}_{r,s'\in D'(\langle s,a\rangle)}[r+\gamma \mathbf{v}(s')] = \mathbb{E}_{\substack{r\sim\mathcal{R}_D(\cdot|\langle s,a\rangle)\\s'\sim P_D(\cdot|\langle s,a\rangle)}}[r+\gamma \mathbf{v}(s')] = [\mathbf{r}_D+\gamma P_D\mathbf{v}](\langle s,a\rangle).$$

Hoeffding's inequality indicates that the mean of bounded random variables will approximate their expected values with high probability. By applying Hoeffding's inequality to each element in the $|S \times A|$ state-action space and employing a union bound, we establish that with probability at least $1 - \delta$,

$$|(\mathbf{r}_D + \gamma P_D \mathbf{v}) - (\mathbf{r}_{D'} + \gamma P_{D'} \mathbf{v})| \le \frac{1}{1 - \gamma} \sqrt{\frac{1}{2} \ln \frac{2|\mathcal{S} \times \mathcal{A}|}{\delta}} \ddot{\mathbf{n}}_{D'}^{-1}$$

We can left-multiply A^{π} and rearrange to get:

$$|A^{\pi}(\mathbf{r}_{D} + \gamma P_{D}\mathbf{v}) - A^{\pi}(\mathbf{r}_{D'} + \gamma P_{D'}\mathbf{v})| \leq \left(\frac{1}{1-\gamma}\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right)A^{\pi}\ddot{\mathbf{n}}_{D'}^{-\frac{1}{2}}$$

then we left-multiply the discounted visitation of π :

$$|(I - \gamma A^{\pi} P_{D'})^{-1} [A^{\pi} (\mathbf{r}_{D} + \gamma P_{D} \mathbf{v}) - A^{\pi} (\mathbf{r}_{D'} + \gamma P_{D'} \mathbf{v})]| \le \left(\frac{1}{1 - \gamma} \sqrt{\frac{1}{2} \ln \frac{2|\mathcal{S} \times \mathcal{A}|}{\delta}}\right) (I - \gamma A^{\pi} P_{D'})^{-1} A^{\pi} \ddot{\mathbf{n}}_{D'}^{-\frac{1}{2}}$$

•

This matrix: $(I - \gamma A^{\pi} P_{D'})^{-1} A^{\pi} \ddot{\mathbf{n}}_{D'}^{-\frac{1}{2}}$, belongs to the space $\mathbb{R}^{|S|}$. By indexing at state *s*, we obtain:

$$(I - \gamma A^{\pi} P_{D'})^{-1} A^{\pi} \ddot{\mathbf{n}}_{D'}^{-\frac{1}{2}}(s) = (1 - \gamma) \sum_{s'} \nu(s' | s_0 = s) \frac{1}{\sqrt{N_{D'}(\langle s, \mu(s) \rangle)}}$$

Finally, by integrate these terms together we have the bound on Lemma 2:

$$|\mathbf{q}_{D'}^*(s,\mu(s)) - \mathbf{q}(s,\mu(s))| \le \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right)\sum_{s'}\nu(s'|s_0=s)\frac{1}{\sqrt{N_{D'}(\langle s',\mu(s')\rangle)}}$$

Given that this inequality is universally applicable to any q, and acknowledging that the heuristic term
h supplied by the LLM serves as a constant within the temporal-difference (TD) update mechanism
of the Q-function, it follows that:

$$|\mathbf{q}_{D'}^{*}(s,\mu(s)) - \hat{\mathbf{q}}(s,\mu(s))| = |\mathbf{q}_{D'}^{*}(s,\mu(s)) - \mathbf{q}(s,\mu(s)) - \mathbf{h}(s,\mu(s))|$$

$$\leq \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right)\sum_{s'}\nu(s'|s_0=s)\frac{1}{\sqrt{N_{D'}(\langle s',\mu(s')\rangle)}}$$

918 A.4 PROOF OF THEOREM 2

⁹²⁰ To get the sample complexity of convergence. By Lemma 2, we have:

$$\begin{aligned} |\mathbf{q}_{D'}^*(s,\mu(s)) - \hat{\mathbf{q}}(s,\mu(s))| &\leq \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right) \sum_{s'} \nu(s'|s_0 = s) \frac{1}{\sqrt{N_{D'}(\langle s',\mu(s')\rangle)}} \\ &= \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right) \sum_{s'} \sqrt{\nu(s'|s_0 = s)} \frac{\sqrt{\nu(s'|s_0 = s)}}{\sqrt{nd_{D'}(s,a)}} \\ &\quad (d_{D'}(s,a) = \frac{N_{D'}(\langle s,a\rangle)}{|D'|}) \\ &= \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right) \sum_{s'} \sqrt{d_{D'}(s,\mu(s))} \frac{\sqrt{d_{D'}(s,\mu(s))}}{\sqrt{nd_{D'}(s,a)}} \\ &\quad (\nu(s)\pi(\mu(s)|s) \approx d_{D'}(s,\mu(s))) \right) \\ &\leq \frac{\left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right)}{\sqrt{n}} \sum_{s'} \sqrt{d_{D'}(s',\mu(s))} \\ &\leq \left(\sqrt{\frac{1}{2}\ln\frac{2|\mathcal{S}\times\mathcal{A}|}{\delta}}\right) \frac{|S|}{\sqrt{n}} \end{aligned}$$

Since D' is sampled iid from replay buffer D, Then, when $n > \mathcal{O}\left(\frac{|S|^2}{2\epsilon^2} \ln \frac{2|S \times A|}{\delta}\right)$, we have $|\mathbf{q}_D^*(s, \mu(s)) - \mathbf{q}^*(s, \mu(s))| \le \epsilon$.

B EXPERIMENT DETAILS

B.1 Q-SHAPING DETAILS

In our experiments, we utilized "gpt-4o" as the language model to provide heuristic Q-values, thereby accelerating the exploration process in the **LLM-TD3** algorithm. The experiments were conducted on a host equipped with a 48-core CPU, 24 GB of GPU memory, and 120 GB of RAM. For complex tasks, the agent took approximately 2 to 4 hours to converge, whereas for simpler tasks, convergence was achieved within 10 to 30 minutes. Table 4 provides a detailed description of the experimental environment.

Table 4: I	Experimental	Environment
------------	--------------	-------------

Resource	Specification
CPU	48-core Intel Xeon E5-2666 v4
GPU	NVIDIA GeForce RTX 4090 (24 GB)
RAM	118.1 GB
Convergence Time (Complex Tasks)	2-4 hours
Convergence Time (Simple Tasks)	10-30 minutes

Hyperparameters LLM-TD3 is built on top of TD3, and doesn't require parameter tuning. In the

baseline implementation, TD3's hyperparameters are also fixed for comparison. The hyperparameters

of LLM-TD3 are detailed in Table 5. Table 6 displays the convergence line for each environment.

72		Table 5: Hyperparameters of LLM-TD3					
7.7				T 7 4			
75	Hy	yperparameter		Value			
76	LI	LM Type		gpt-40			
77	St	art Timesteps		5000			
70	Ev	aluation Frequency	:	5,000			
/8	Ex	ploration Noise (St	d) (0.1			
79	Ba	atch Size		256			
80	Di	scount Factor γ		0.99			
31	Ta	rget Network Updat	te Rate (Tau)	0.005			
32	Po	olicy Noise		0.2			
3	No	bise Clip		0.5			
4	Pc	licy Update Freque	ncy	2			
5	λ_1	λ_2		100,10	`		
6	Hi	aden Layer Size		512 (10,240 for Humanoid)		
7							
8		Table 6: Conve	ergence Line for	r Each Environment			
9			<i>6 · · · · · · · · · · · · · · · · · · ·</i>				
0		Environment		Convergence Line			
1							
2		Ant-v4		4480			
3		HalfCheetah-v4		8800			
4		Hopper-v4		2560			
5		Humanoid-v4	4	4000			
6		InvertedPendulu	m-v4	800			
7		Pendulum-v1		-200			
0		Walker2D-v4	ntinuous	5700			
5 D		Drower Open To	nunuous	3200			
9		Window-Close-	15KI Fack1	3200			
JU		Button_Press_Ta	ek1	3200			
01		Sween-Into-Task	71	2800			
)2		Door-Close-Tasl	<1 <1	3200			
)3		Handle-Press-Ta	sk1	3200			
)4		Basketball-V2-T	ask1	360			
)5		Coffee-Button-V	2-Task1	2960			
)6		Soccer-V2-Task	1	1600			
)7		PvFlvt/OuadX-F	Sall-In-Cup-V2	3840			
8		PyFlyt/OuadX-F	Pole-Balance-V2	2 1600			
9		PyFlyt/QuadX-H	Hover-V2	880			
0							
1							
2							
13	B.2 BASELINE	DETAILS					
1/1	XX7 . 11 7 . 1	1		1			
15	We use table 7 to 1	st the open source r	epositories of th	ie algorithms used in the e	xperiment, Figu		
10	8 to present the hyp	perparameters of cle	anRL_SAC, and	d Figure 9 to present the h	yperparameters		
10	cleanKL_PPO.						
17							
8		Table	e 7: Baseline Co	ode Source			
19							
1U 24		Algorithm	Code Reposito	ory			
. 1		cleanRL PPO	https://github.	com/vwxyzjn/cleanrl			
2		TD3	https://github.	com/sfujim/TD3			

Table 5: Hyperparameters of LLM-TD3

https://github.com/sfujim/TD3

https://github.com/vwxyzjn/cleanrl

DDPG

cleanRL_SAC

1023

1024

Table 8: Hyperparameter	rs of SAC	Table 9: Key Hyperparameter	Table 9: Key Hyperparameters of PPC		
Hyperparameter	Value	Hyperparameter	Value		
Critic Learning Rate	3e-3	Learning Rate	3e-4		
Actor Learning Rate	3e-4	Num Steps	2048		
Entropy Target	$-\dim(\mathcal{A})$	Total Timesteps	1e6		
Policy Update Frequency	1	Gamma (Discount Factor)	0.99		
Reward Scale	$\frac{1}{2}$, 1	GAE Lambda	0.95		
Hidden Layer Size	128	Clip Coefficient	0.2		

1038 B.3 DETAILS OF IMPLEMENTING LLM-BASED REWARD SHAPING METHODS

In this experiment, we evaluate Q-shaping against Text2Reward (T2R) (Xie et al.) and Eureka (Ma et al., 2023) to compare LLM-based reward shaping approaches.

Text2Reward (T2R): Text2Reward is a framework designed to address the challenge of reward shaping in reinforcement learning by automating the generation of dense, interpretable reward codes using large language models (LLMs). This method demonstrates effectiveness across various robotic and locomotion tasks, achieving success rates comparable to or exceeding those obtained with expert-designed reward codes (Xie et al.). In our experiment, we implement T2R using the provided prompt available at GitHub link and Soft Actor-Critic (SAC) as the RL backbone. The hyperparameters listed in Table 10 are used for the implementation of SAC in the Text2Reward experiment.

Table 10: Key Hyperparameters for SAC Implementation

Hyperparameter	Value
Batch Size	512
Policy Network Architecture	[256, 256, 256]
Discount Factor (γ)	0.99
Learning Rate	0.0003
Soft Update Coefficient (τ)	0.005
Learning Starts (steps)	25,000
Entropy Coefficient (α)	auto_0.1

Eureka:

- Eureka is a reward design algorithm that leverages the capabilities of LLMs for evolutionary optimization of reward functions. It uses the environment code as context, generating executable reward functions in a zero-shot manner, and iteratively improves them through reflection-based feedback and evolutionary search. Eureka's robust framework has been validated across a wide range of RL tasks, outperforming expert-designed rewards in many scenarios (Ma et al., 2023).
- 2. Eureka is originally designed to operate within the Isaac Gym simulator, adaptations were necessary for our experiments to integrate Eureka's functionality with our environment. Specifically, the prompt for Eureka was tailored into two configurations: one for initial code generation and another for refining the code based on feedback. These prompts are detailed in Eureka Prompt 1: Code Generation and Eureka Prompt 1: Reflection. The first prompt facilitates the generation of foundational reward programs, while the second focuses on optimizing these codes iteratively to align better with experimental objectives.

In our implementation of Eureka, we configured the iterative batch size (K) to 8 and the search iterations (N) to 5. Table 11 summarizes the results of each evolutionary iteration. It shows agent performance at each run and the improvement of each evolution.

1080	Eureka Prompt 1: Code Generation
1081	1 A A A A A A A A A A A A A A A A A A A
1082	You are a reward engineer trying to write reward functions to solve reinforcement learning
1083	tasks as effective as possible. Your goal is to write a reward function for the environment
1084	that will help the agent learn the task described in text. Your reward function should use
1085	useful variables from the environment as inputs. As an example, the reward function
1086	signature can be:
1087	<pre>def compute_reward_shaped(obs: torch.Tensor, action: torch.</pre>
1088	Tensor) ->
1089	3 Tuple[[float, Dict[str. float]]]
1090	4
1000	5
1092	<pre>6 return reward, { }</pre>
1093	the obs shape is {batch_size, obs_dim} and action shape is {batch_size, action_dim}. and
1094	batch_size is 1. Make sure any new tensor or variable you introduce is on the same device
1095	as the input tensors.
1097	The Python environment is {task_obs_code_string}. Write a reward function for the
1098	tollowing task: {task_description}.
1099	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
1100	a python code string: """ ""
1101	Some helpful tips for writing the reward function code:
1102	1. You may find it helpful to normalize the reward to a fixed range by applying
1103	transformations like torch.exp to the overall reward or its components.
1104	2. If you choose to transform a reward component, then you must also introduce a
1105	temperature parameter inside the transformation function; this parameter must
1106	be a named variable in the reward function and it must not be an input variable.
1107	Each transformed reward component should have its own temperature variable.
1108	3. Make sure the type of each input variable is correctly specified; a float input
1109	variable should not be specified as torch. Tensor.
1110	4. Most importantly, the reward code's input variables must contain only attributes
1111	prefix self.) Under no circumstance can you introduce new input variables
1112	prenx Serre). Onder no encumstance can you muoduce new input variables.
1113	
1114	Eureka Prompt 2: Reflection
1110	1
1117	You are a reward engineer trying to write reward functions to solve reinforcement learning
1118	tasks as effective as possible. Your goal is to write a reward function for the environment
1119	that will help the agent learn the task described in text. Your reward function should use
1120	signature can be:
1121	
1122	<pre>1 def compute_reward_shaped(obs: torch.Tensor, action: torch. Tensor) > Turple[fleat_Dist[str_fleat]];</pre>
1123	lensor) -> lupie[lioat, Dict[str, lioat]]:
1124	<pre>return reward, {}</pre>
1125	the one shape is (batch size one dim) and action shape is (batch size action dim) and
1126	hatch size is 1 Make sure any new tensor or variable you introduce is on the same device
1127	as the input tensors.
1128	The Python environment is {task_obs_code_string}. Write a reward function for the
1129	following task: {task_description}.
1130	The output of the reward function should consist of two items: (1) the total reward, (2) a
1131	dictionary of each individual reward component. The code output should be formatted as
1132	a python code string: ""python "".
1133	some neiprar ups for writing the reward runction code.

1104	
1134	1. You may find it helpful to normalize the reward to a fixed range by applying
1136	transformations like torch. exp to the overall reward or its components.
1137	2. If you choose to transform a reward component, then you must also introduce a
1138	temperature parameter inside the transformation function; this parameter must
1139	be a named variable in the reward function and it must not be an input variable.
1140	Each transformed reward component should have its own temperature variable.
1141	3. Make sure the type of each input variable is correctly specified; a float input variable should not be specified as torrch. Tensor
1142	A Most importantly the reward code's input variables must contain only attributes
1143	of the provided environment class definition (namely, variables that have the
1144	prefix self.). Under no circumstance can you introduce new input variables.
1145	{the best code}
1146	We trained a RL policy using the provided reward function code and tracked the values
1147	of the individual components in the reward function as well as global policy metrics such
1148	as success rates and episode lengths after every {epoch_freq} epochs and the maximum,
1149	mean, minimum values encountered:
1150	{data}
1151 1152	that can better solve the task. Some helpful tips for analyzing the policy feedback:
1153	1. If the success rates are always near zero, then you must rewrite the entire reward function
1154	2. If the values for a certain reward component are near identical throughout, then
1156	this means RL is not able to optimize this component as it is written. You may
1157	consider:
1158	(a) Changing its scale or the value of its temperature parameter,
1159	(b) Re-writing the reward component,
1160	(c) Discarding the reward component.
1161	3. If some reward components' magnitude is significantly larger, then you must
1162	re-scale its value to a proper range.
1163	Please analyze each existing reward component in the suggested manner above first, and
1164	then write the reward function code.
1165	
1166	
1167	

Table 11: Average episodic returns for the task *drawer-open* across different evolution rounds (r_x) and agents (a_x) , evaluated at 200k training steps. r_x denotes the evolution round, and a_x represents the agent in that round. The column *best* indicates the best-performing agent in each round.

Round (r_x)	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8]
r_1	1018.77	371.74	1931.21	2117.50	2145.50	2373.34	2483.56	1704.40	
r_2	1964.78	357.45	1704.17	2268.05	1869.65	2073.90	2124.35	1561.83	
r_3	2163.36	1123.12	801.96	1993.93	2163.62	1904.10	850.56	428.80	
r_4	2142.97	839.31	1260.82	1445.53	1665.38	1470.83	433.12	1063.76	
r_5	928.98	1392.44	1761.59	2123.26	2308.81	1348.77	698.31	1888.74	

¹¹⁷⁸

1179

1180 1181

1182

C PROMPT DETAILS FOR THE Q-SHAPING FRAMEWORK

The Q-shaping framework necessitates a general template to guide the code generation provided by
large language models (LLMs). This template requires three key components: (1) the code template,
(2) the environment description, and (3) the environment configuration file.

Below is a comprehensive overview of the general template:

1188	General Prompt
1190	Now word to prove a give of and hand on the description of the environment on the
1191	You need to generate a piece of code based on the description of the environment of the
1192	
1193	The purpose of this code is to provide a suitable O value for (s, a) that you consider good
1194	based on the information provided. For bad (s, a), you can assign a Q-value of 0 or a
1195	lower value to discourage the robot from taking this action.
1196	Requirements:
1197	1. In short, your task is to convert the task description into a Python-style $Q(s, a)$
1198	2. The environment description typically provides the obs_dim and action_dim, along
1199	with the conditions for terminal states and truncation. Your task is to penalize behaviors
1200	3 If you are confident, you can use your knowledge to generate (s a O) values that you
1201	believe may lead to success or failure. 4. The code returns s. a. g. targets
1202	5. Generate two functions, def good_Q(self, batch_size), def bad_Q(self, batch_size)
1203	6.TIPS: Action is more important than state, so you should focus on encouraging actions
1204	that lead to success and discouraging actions that lead to failure.
1205	7. When designing bad Q-values, there are no bad states, only bad actions. You need to
1206	clearly identify which state-action pairs lead to termination and avoid those actions.
1207	8. If the description mentions states that lead to termination, you should include them in the had Ω values as assigning a Ω value of 0 to termination states usually associated
1208	learning
1209	9 You can try to encourage as many (s a) pairs as possible to guide the agent to explore
1210	directions that you believe will lead to success.
1211	10. You should provide a complete class definition, including theinit, goodQ, and
1212	badQ methods, without omitting any of them.
1213	
1214	{code template}
1215	(in the second data)
1216	{environment description}
1217	{environment config file}
1218	
1219	
1220	
1661	

C.1 ILLUSTRATIVE EXAMPLE: Q-SHAPING FRAMEWORK IN ACTION

To provide a concrete understanding of the Q-shaping framework, we present an example using the robotic arm task "handle-press-v2". This example illustrates the application of the general template outlined earlier and demonstrates how the three key components—code template, environment description, and environment configuration file—come together to generate (s,a,Q) pairs that effectively guide agent behavior.

1229 1230

1222

1223 1224

1225

1226

1227

1228

1231 C.1.1 ENVIRONMENT DESCRIPTION

The Meta-World benchmark is a suite of 50 diverse robotic manipulation tasks designed to evaluate reinforcement learning (RL) and meta-reinforcement learning (meta-RL) algorithms. In Yu et al. (2020), the authors introduce a simulated Sawyer robotic arm and provide detailed definitions of the observation space, action space, and evaluation metrics.

For the purpose of this paper, we focus on Section 4.1 Actions, Observations, and Rewards from Yu et al. (2020), which outlines the design of the state space, action space, and reward functions. These details are critical for understanding how to guide large language models (LLMs) to generate high-quality (s, a, Q) pairs.

Environment Description

1242

1243

1282

1244	4.1 Actions Observations and Rewards. In order to represent policies for multiple
1245	tasks with one model, the observation and action snaces must contain significant shared
1246	structure across tasks. All of our tasks are performed by a simulated Sawyer robot.
1247	The action space is a 2-tuple consisting of the change in 3D space of the end-effector
1248	followed by a normalized torque that the gripper fingers should apply. The actions in
1249	this space range between -1 and 1. For all tasks, the robot must either manipulate one
1250	object with a variable goal position, or manipulate two objects with a fixed goal position.
1251	The observation space is represented as a 6-tuple of the 3D Cartesian positions of the
1252	end-effector, a normalized measurement of how open the gripper is, the 3D position of
1253	the first object, the quaternion of the first object, the 3D position of the second object, the
1254	quaternion of the second object, all of the goal. If there is no second object or the goal is not
1255	meant to be included in the observation, then the quantities corresponding to them are zeroed out. The observation space is always 39 dimensional.
1256	
1257	
1258	Designing reward functions for Meta-World requires two major considerations. First, to
1259	guarantee that our tasks are within the reach of current single-task reinforcement learning
1260	algorithms, which is a prerequisite for evaluating multi-task and meta-RL algorithms,
1261	we design well-shaped reward functions for each task that make each of the tasks at least
1262	individually solvable.
1263	Man importantly the proved functions must subility shared structure serves tasks
1264	More importantly, the reward functions must exhibit shared structure across tasks.
1265	varying reward scales or structures can make the tasks appear completely distinct for the
1266	learning algorithm, masking their shared structure and leading to preferences for tasks
1267	with high-magnitude rewards.
1268	
1269	Accordingly, we adopt a structured, multi-component reward function for all tasks, which
1270	leads to effective policy learning for each of the task components. For instance, in a task
1271	that involves a combination of reaching, grasping, and placing an object, let $o \in \mathbb{R}^3$ be
1272	the object position, where $o = (o_x, o_y, o_z)$, $h \in \mathbb{R}^3$ be the position of the robot's gripper,
1273	$z_{\text{target}} \in \mathbb{R}$ be the target height of lifting the object, and $g \in \mathbb{R}^3$ be goal position. With
1274	the above definition, the multi-component reward function R is the combination of a
1275	reaching reward, a grasping reward, and a placing reward or subsets thereof for simpler
1276	across all tasks have a similar magnitude that ranges between 0 and 10, where 10 always
1277	corresponds to the reward-function being solved and conform to similar structure as
1278	desired. The full form of the reward function and a list of all task rewards is provided in
1279	Appendix.
1280	
1281	

C.1.2 ENVIRONMENT CONFIGURATION FILE

The primary purpose of the configuration file is to specify the target object's location and the initial position of the robotic arm's gripper. This information can assist the LLM in generating movement direction vectors that lead to effective actions.

```
1 from ____future___ import annotations
1287
1288
    3 from typing import Any
1289
    4
1290 5 import numpy as np
1291 6 import numpy.typing as npt
    7 from gymnasium.spaces import Box
1292
    8
1293
    9 from metaworld.envs.asset_path_utils import full_v2_path_for
1294 10 from metaworld.envs.mujoco.sawyer_xyz.sawyer_xyz_env import RenderMode,
1295
         SawyerXYZEnv
    11 from metaworld.envs.mujoco.utils import reward_utils
```

```
1296
     12 from metaworld.types import InitConfigDict
1297
    13
1298 14
1299 15 class SawyerHandlePressEnvV2 (SawyerXYZEnv):
            TARGET_RADIUS: float = 0.02
1300 <sup>16</sup>
1301 <sup>17</sup>
           def __init__(
     18
1302
               self,
     19
1303 20
                render_mode: RenderMode | None = None,
1304 21
                camera_name: str | None = None,
1305 <sup>22</sup>
                camera_id: int | None = None,
           ) -> None:
1306<sup>23</sup>
                hand_low = (-0.5, 0.40, 0.05)
     24
1307
     25
               hand_high = (0.5, 1.0, 0.5)
1308 26
                obj_low = (-0.1, 0.8, -0.001)
                obj_high = (0.1, 0.9, 0.001)
1309 27
                goal_low = (-0.1, 0.55, 0.04)
1310 <sup>28</sup>
                goal_high = (0.1, 0.70, 0.08)
1311 <sup>29</sup>
     30
1312 31
                super().__init__(
1313 32
                     hand_low=hand_low,
1314 33
                     hand_high=hand_high,
1315 <sup>34</sup>
                     render_mode=render_mode,
1316 <sup>35</sup>
                     camera_name=camera_name,
                     camera_id=camera_id,
     36
1317 37
                )
1318 38
1319 39
                self.init_config: InitConfigDict = {
                     "obj_init_pos": np.array([0, 0.9, 0.0]),
1320 40
                     "hand_init_pos": np.array(
1321 <sup>41</sup>
                          (0, 0.6, 0.2),
     42
1322 43
                     ),
1323 44
                }
                self.goal = np.array([0, 0.8, 0.14])
1324 45
                self.obj_init_pos = self.init_config["obj_init_pos"]
1325 <sup>46</sup>
1326 47
                self.hand_init_pos = self.init_config["hand_init_pos"]
     48
1327 49
                self._random_reset_space = Box(
1328 50
                     np.array(obj_low), np.array(obj_high), dtype=np.float64
1329 51
                )
                self.goal_space = Box(np.array(goal_low), np.array(goal_high),
1330 52
           dtype=np.float64)
1331
     53
1332 54
            @property
1333 55
            def model_name(self) -> str:
                return full_v2_path_for("sawyer_xyz/sawyer_handle_press.xml")
1334 56
1335 57
1336 58
            @SawyerXYZEnv._Decorators.assert_task_is_set
     59
            def evaluate_state(
1337
                self, obs: npt.NDArray[np.float64], action: npt.NDArray[np.
     60
1338
            float32]
            ) -> tuple[float, dict[str, Any]]:
1339 61
                 (
1340 <sup>62</sup>
1341 <sup>63</sup>
                     reward,
                     tcp_to_obj,
     64
1342 65
                     _/
1343 66
                     target_to_obj,
1344 67
                     object_grasped,
1345 <sup>68</sup>
                    in_place,
     69
                ) = self.compute_reward(action, obs)
1346
     70
1347 71
                info = \{
1348 72
                     "success": float(target_to_obj <= self.TARGET_RADIUS),
                     "near_object": float(tcp_to_obj <= 0.05),</pre>
1349 73
                     "grasp_success": 1.0,
     74
```

```
1350
                     "grasp_reward": object_grasped,
     75
1351
                     "in_place_reward": in_place,
     76
1352 77
                     "obj_to_target": target_to_obj,
1353 78
                     "unscaled_reward": reward,
1354 <sup>79</sup>
                }
1355 <sup>80</sup>
                return reward, info
     81
1356
     82
1357 83
            0property
1358 84
            def _target_site_config(self) -> list[tuple[str, npt.NDArray[Any]]]:
1359 <sup>85</sup>
                 return []
1360 <sup>86</sup>
            def _get_pos_objects(self) -> npt.NDArray[Any]:
     87
1361
                return self._get_site_pos("handleStart")
     88
1362 89
            def _get_quat_objects(self) -> npt.NDArray[Any]:
1363 90
                return np.zeros(4)
1364 91
1365 92
     93
            def _set_obj_xyz(self, pos: npt.NDArray[Any]) -> None:
1366 94
                qpos = self.data.qpos.flat.copy()
1367 95
                qvel = self.data.qvel.flat.copy()
                qpos[9] = pos
1368 96
1369 97
                qvel[9] = 0
1370 98
                self.set_state(qpos, qvel)
    99
1371 100
            def reset_model(self) -> npt.NDArray[np.float64]:
1372 101
                self._reset_hand()
1373 102
1374 <sup>103</sup>
                self.obj_init_pos = self._get_state_rand_vec()
1375<sup>104</sup>
                self.model.body("box").pos = self.obj_init_pos
                self._set_obj_xyz(np.array(-0.001))
    105
1376 106
                self._target_pos = self._get_site_pos("goalPress")
1377 107
                self.maxDist = np.abs(
                     self.data.site("handleStart").xpos[-1] - self._target_pos[-1]
1378 108
                )
1379 <sup>109</sup>
1380<sup>110</sup>
                self.target_reward = 1000 * self.maxDist + 1000 * 2
    111
                self._handle_init_pos = self._get_pos_objects()
1381 112
1382 113
                return self._get_obs()
1383 114
1384 <sup>115</sup>
            def compute_reward(
1385<sup>116</sup>
                self, actions: npt.NDArray[Any], obs: npt.NDArray[np.float64]
            ) -> tuple[float, float, float, float, float, float]:
1386 118
                assert (
1387 119
                     self._target_pos is not None
                ), "'reset_model()' must be called before 'compute_reward()'."
1388 120
                del actions
1389 <sup>121</sup>
1390<sup>122</sup>
                obj = self._get_pos_objects()
                tcp = self.tcp_center
    123
1391 124
                target = self._target_pos.copy()
1392 125
                target_to_obj = obj[2] - target[2]
1393 126
1394 <sup>127</sup>
                target_to_obj = np.linalg.norm(target_to_obj)
1395<sup>128</sup>
                target_to_obj_init = self._handle_init_pos[2] - target[2]
                target_to_obj_init = np.linalg.norm(target_to_obj_init)
    129
1396 130
1397 131
                in_place = reward_utils.tolerance(
1398 132
                     target_to_obj,
1399 <sup>133</sup>
                     bounds=(0, self.TARGET_RADIUS),
1400<sup>134</sup>
                     margin=abs(target_to_obj_init - self.TARGET_RADIUS),
                     sigmoid="long_tail",
1401 136
                )
1402 137
1403 138
                handle_radius = 0.02
                tcp_to_obj = float(np.linalg.norm(obj - tcp))
    139
```

```
1404
                 tcp_to_obj_init = np.linalq.norm(self._handle_init_pos - self.
    140
1405
            init_tcp)
1406 141
               reach = reward_utils.tolerance(
                     tcp_to_obj,
1407 142
1408<sup>143</sup>
                     bounds=(0, handle_radius),
1409<sup>144</sup>
                     margin=abs(tcp_to_obj_init - handle_radius),
                     sigmoid="long_tail",
1410 146
                )
1411 147
                 tcp_opened = 0
1412 148
                 object_grasped = reach
1413 <sup>149</sup>
1414<sup>150</sup>
                reward = reward_utils.hamacher_product(reach, in_place)
    151
                reward = 1.0 if target_to_obj <= self.TARGET_RADIUS else reward</pre>
1415 152
                reward *= 10
1416 153
                 return (reward, tcp_to_obj, tcp_opened, target_to_obj,
           object_grasped, in_place)
1417
                             Listing 1: Config file for sawyer-handle-press-v2
1418
```

1421 C.1.3 CODE TEMPLATE

1419 1420

In the code template section, a Python-style code snippet and its explanation are provided. The Pythonstyle code defines the expected output format of the LLM, while the accompanying explanation helps
the LLM better understand the structure and purpose of the code.

```
1425
            Code Template
1426
1427
            For example:
1428
            In the DrawerOpen environment, actions are encouraged when they involve moving the
1429
            gripper towards the handle and closing the gripper. Specifically:
1430
1431
            Encouraged actions (good Q):
1432
1433
            Movement in the direction of the handle (positive y-direction).
1434
            Closing the gripper, especially when the gripper is close to the handle.
            Discouraged actions (bad Q):
1435
1436
            Movement away from the handle (negative y-direction).
1437
            Opening the gripper when it is near the handle, or further opening it when it's already
1438
            open.
1439
1440
          1 class DrawerOpen:
1441
                def ___init___(self):
          2
1442
                      self.obs_dim = 39 # Observation space dimension
          3
1443
                     self.action_dim = 4 # Action space dimension (dx, dy,
          4
1444
                dz, gripper torque)
1445
                      self.maxDist = 0.2 # Maximum distance for drawer
                opening
1446
                     self.target_reward = 1000 * self.maxDist + 1000 * 2
          6
1447
                     self.close_gripper_threshold = 0.05 # Distance
           7
1448
                threshold to encourage closing the gripper
1449
          8
1450
                 def good_Q(self, batch_size):
          9
                     actions = []
          10
1451
                      states = []
          11
1452
                      q_targets = []
1453
          13
                     for _ in range(batch_size):
1454
                          # Generate a state where the gripper is approaching
          14
1455
                the handle
          15
                          handle_{pos} = np.array([0.0, 0.74, 0.09])
1456
                 Approximate handle position
1457
```

1458	
1459	16 # Start gripper at a position slightly away from the
1460	handle
1461	<pre>17 gripper_pos = handle_pos + np.random.uniform(-0.15, 0.15 cigo=2)</pre>
1462	gripper open = np.random.uniform(0.0, 0.5) #
1463	Gripper partially closed
1464	19
1465	20 # Construct the observation
1466	21 obs = np.zeros(self.obs_dim) abs[:3] = gripper pos_ # Cripper position
1467	22 obs[3] = gripper_pos # Gripper position 23 obs[3] = gripper open # Gripper state
1468	24 obs[4:7] = handle_pos # Handle position
1469	<pre>25 obs[7:] = np.random.uniform(-0.1, 0.1, size=self.</pre>
1470	obs_dim - 7) # Other observations
1471	26 # Concrate actions that move the gripper towards the
1472	handle (positive v movement)
1473	<pre>28 direction_to_handle = handle_pos - gripper_pos</pre>
1474	<pre>29 distance_to_handle = np.linalg.norm(</pre>
1475	direction_to_handle)
1476	action direction = direction to handle /
1477	distance_to_handle
1478	32 else:
1479	<pre>33 action_direction = np.zeros(3)</pre>
1480	action_magnitude = np.random.uniform(0.05, 0.1)
1481	action magnitude
1482	36
1483	37 # Encourage closing the gripper when close to the
1484	handle
1485	<pre>38 II distance_to_nandle < sell.close_gripper_threshold .</pre>
1486	<pre>39 gripper_action = np.random.uniform(0.5, 1.0) #</pre>
1487	Close the gripper more aggressively
1488	40 else:
1409	Keep the gripper partially open
1490	42
1431	43 action = np.concatenate((
1/03	44 action_movement, # Move towards the handle
1404	45 [gripper_action] # Gripper action
1495	47
1496	48 # Calculate a higher Q-value for actions that reduce
1497	the distance to the handle and close the gripper
1498	49 new_gripper_pos = gripper_pos + action[:3] 50 new_distance_to_handle_=_np_lipplg_norm(handle_pos
1499	new gripper pos)
1500	if new_distance_to_handle < self.
1501	<pre>close_gripper_threshold and gripper_action > 0.5:</pre>
1502	52 q_value = (1.0 - new_distance_to_handle / self.
1503	maxDist) * 15.0 # Higner reward for closing hear the handle
1504	g_value = max(0.0, 1.0 - new_distance_to_handle
1505	/ self.maxDist) * 10.0
1506	55
1507	56 states.append(obs)
1508	58 g targets.append(g value)
1509	59
1510	60 # Convert lists to tensors
1511	61 states = torch.tensor(states, dtype=torch.float32)
	1

actions = torch.tensor(actions, dtype=torch.float32). view(batch_size, self.action_dim) q_targets = torch.tensor(q_targets, dtype=torch.float32) .view(-1, 1) return states, actions, q_targets def bad_Q(self, batch_size): actions = [] states = [] q_targets = [] for _ in range(batch_size): # Generate a state where the gripper is far from the handle gripper_pos = np.array([0.0, 0.5, 0.2]) + np.random. uniform(-0.1, 0.1, size=3) gripper_open = np.random.uniform(0.5, 1.0) # Gripper open handle_pos = np.array([0.0, 0.74, 0.09]) # Handle position remains the same # Construct the observation obs = np.zeros(self.obs_dim) obs[:3] = gripper_pos # Gripper position obs[3] = gripper_open # Gripper state obs[4:7] = handle_pos # Handle position obs[7:] = np.random.uniform(-1.0, 1.0, size=self. obs_dim - 7) # Other observations # Generate actions that move the gripper away from the handle or open the gripper further # Positive y direction moves closer to the handle, so negative y moves away and should be punished action = np.array([np.random.uniform(-1.0, 1.0), # Random x movement np.random.uniform(-1.0, 0.0), # Negative y movement (away from handle) np.random.uniform(-1.0, 1.0), # Random z movement np.random.uniform(-1.0, 0.0) # Torque to open gripper 1) # Assign a low Q-value for bad actions (moving away from the handle) $q_value = 0.0$ states.append(obs) actions.append(action) q_targets.append(q_value) # Convert lists to tensors states = torch.tensor(states, dtype=torch.float32) actions = torch.tensor(actions, dtype=torch.float32). view(batch_size, self.action_dim) q_targets = torch.tensor(q_targets, dtype=torch.float32) .view(-1, 1) return states, actions, q_targets Listing 2: Code template

1566 In the example above, the good_Q method assigns high Q-values to actions likely to lead to success, 1567 whereas bad_Q discourages actions leading to termination. This approach aligns with the principles 1568 outlined in the general template. 1569

1570 C.1.4 LLM OUTPUT

1571

1577

Based on the requirements and suggestions outlined in the General Template, as well as 1572 the observation and action space definitions provided in the Environment Description and 1573 Configuration File, and the Code Template discussed above, the LLM can generate good 1574 Q pairs and bad Q pairs that meet the specified requirements. 1575

Below is the generated code: 1576

```
1 class HandlePress:
1578
     2
           def __init__(self):
               self.obs_dim = 39 # Observation space dimension
1579
     3
               self.action_dim = 4 # Action space dimension (dx, dy, dz,
1580
     4
          gripper torque)
1581
               self.TARGET_RADIUS = 0.02 # Target radius for pressing the
1582
          handle
1583
               self.handle_init_pos = np.array([0.0, 0.9, 0.0]) # Approximate
    6
          handle start position
1584
               self.target_pos = np.array([0.0, 0.8, 0.14]) # Target position
     7
1585
           after pressing
1586
               self.max_dist = np.abs(self.handle_init_pos[2] - self.target_pos
1587
           [21)
1588
     9
               self.gripper_init_pos = np.array([0.0, 0.6, 0.2]) # Gripper
           initial position
1590 <sup>10</sup>
           def good_Q(self, batch_size):
    11
1591
               actions = []
    12
1592
               states = []
1593 14
               q_targets = []
1594 15
               for _ in range(batch_size):
                    # Generate gripper position near the handle
1595 <sup>16</sup>
                   gripper_pos = self.handle_init_pos + np.random.uniform(-0.8,
1596<sup>17</sup>
           0.8, size=3)
1597
    18
                    # gripper_pos[2] += 0.1 # Ensure gripper is above the handle
1598
            in z-axis
                   gripper_open = np.random.uniform(-1, 1) # Gripper partially
1599 19
           closed
1600
    20
1601
                    # Construct the observation
1602
                   obs = np.zeros(self.obs_dim)
    22
1603
                    obs[:3] = gripper_pos # Gripper position
                    obs[3] = gripper_open # Gripper state
1604 24
                    obs[4:7] = self.handle_init_pos # Handle position
1605 <sup>25</sup>
                   obs[7:] = np.random.uniform(-0.3, 0.3, size=self.obs_dim - 7)
1606 <sup>26</sup>
             # Other observations
1607
    27
1608 28
                    # Generate actions that move the gripper towards the handle
           and press it down
1609
                   x_move = np.random.uniform(-0.001, 0.001)
1610 <sup>29</sup>
                   y_move = np.random.uniform(-0.001, 0.001)
    30
1611
                    z_move = np.random.uniform(-0.1, -0.08)
    31
1612
                   action_movement = np.array([x_move, y_move, z_move])
    32
1613 33
                    # Apply torque to press the handle
1614 34
                   gripper_action = np.random.uniform(-0.001, 0.001) # Apply
           torque to press
1615
    35
1616
                    action = np.concatenate((
    36
1617
    37
                        action_movement, # Move towards the handle and press
1618
           down
1619 38
                        [gripper_action] # Gripper torque action
    39
                   ))
```

```
1620
     40
1621
                    # Calculate Q-value based on how much the handle is pressed
    41
1622
           towards the target
1623 42
                    handle_movement = action_movement[2] # Z-axis movement (
           pressing down)
1624
                    new_handle_pos_z = self.handle_init_pos[2] + handle_movement
    43
1625
                    handle_to_target = np.abs(new_handle_pos_z - self.target_pos
    44
1626
           [2])
1627 45
1628 46
                    q_value = 10
1629 47
                    states.append(obs)
                    actions.append(action)
    48
1630
                    q_targets.append(q_value)
    49
1631
    50
1632 51
                # Convert lists to tensors
                states = torch.tensor(states, dtype=torch.float32)
1633 52
                actions = torch.tensor(actions, dtype=torch.float32).view(
1634 53
           batch_size, self.action_dim)
1635
    54
                q_targets = torch.tensor(q_targets, dtype=torch.float32).view(-1,
1636
            1)
1637 <sub>55</sub>
1638 56
                return states, actions, q_targets
1639 57
1640 58
           def bad_Q(self, batch_size):
                actions = []
    59
1641
                states = []
    60
1642 61
                q_targets = []
1643 62
                for _ in range(batch_size):
                    # Generate gripper position far from the handle
1644 <sup>63</sup>
                    gripper_pos = self.gripper_init_pos + np.random.uniform(-0.8,
1645 <sup>64</sup>
            0.8, size=3)
1646 65
                    gripper_open = np.random.uniform(-1, 1) # Gripper open
1647 66
1648 67
                    # Construct the observation
                    obs = np.zeros(self.obs_dim)
1649 <sup>68</sup>
                    obs[:3] = gripper_pos # Gripper position
    69
1650
                    obs[3] = gripper_open # Gripper state
    70
1651
                    obs[4:7] = self.handle_init_pos # Handle position
1652 72
                    obs[7:] = np.random.uniform(-1.0, 1.0, size=self.obs_dim - 7)
             # Other observations
1653
1654 <sup>73</sup>
                    # Generate actions that move the gripper away from the handle
    74
1655
            or lift it up
1656 75
                    action = np.array([
1657 76
                         np.random.uniform(-1.0, 1.0), # Random x movement
                         np.random.uniform(0.8, 1.0), # Random y movement
1658 77
                         np.random.uniform(0.2, 1.0), # Positive z movement (move
1659 <sup>78</sup>
            up)
1660
                         np.random.uniform(-1.0, 1) # Negative torque (release
    79
1661
           pressure)
1662 80
                    1)
1663 81
                    # Assign a low Q-value for bad actions
1664 <sup>82</sup>
                    q_value = -10
    83
1665
    84
1666
                    states.append(obs)
    85
1667 86
                    actions.append(action)
1668 87
                    q_targets.append(q_value)
1669 <sup>88</sup>
                # Convert lists to tensors
    89
1670
                states = torch.tensor(states, dtype=torch.float32)
    90
1671 91
                actions = torch.tensor(actions, dtype=torch.float32).view(
1672
           batch_size, self.action_dim)
1673 92
                q_targets = torch.tensor(q_targets, dtype=torch.float32).view(-1,
            1)
```

return states, actions, q_targets Listing 3: LLM generated code for HandlePress LIMITATIONS AND DISCUSSION D In this section, we discuss several limitations of our work and propose promising directions for future research to address these challenges. One limitation of the Q-shaping framework is its dependence on the understanding capabilities of large language models (LLMs). For tasks where the LLM cannot fully comprehend the dynamics of movement or the control-specific requirements of the task, the framework fails to deliver meaningful improvement. This limitation constrains the applicability of Q-shaping to domains where task requirements can be effectively interpreted by the LLM. Models with stronger reasoning capabilities may be needed to generate valid state-action pairs. Another limitation is the difficulty of scaling Q-shaping to visual or real-world settings. The frame-work requires models capable of generating states, but current technology lacks models that can simultaneously process textual and visual inputs and output comprehensive state-action descriptions. This gap restricts the ability of Q-shaping to operate effectively in environments where visual data is a critical component. Future progress in multimodal modeling, such as vision-language models that integrate text and images, could alleviate this challenge by enabling richer state representations. By addressing these limitations, Q-shaping has the potential to evolve into a more versatile framework capable of operating across diverse tasks and environments, ultimately advancing its impact on reinforcement learning research.