

Reward-Zero: Language Embedding Driven Implicit Reward Mechanisms for Reinforcement Learning

Anonymous authors

Paper under double-blind review

Keywords: Sparse reward, Language-conditioned reward, Auxiliary reward shaping, Efficient learning

Summary

Reinforcement learning agents often struggle with sparse or poorly shaped reward signals, and hand-crafting dense rewards for each new task is labor-intensive and error-prone. We introduce Reward-Zero, an implicit reward mechanism that derives dense progress signals from natural-language goal descriptions using pre-trained vision-language embeddings. Given only a textual goal (e.g., “The cabinet drawer is fully open”) and raw visual observations, Reward-Zero computes a potential function via CLIP image–text similarity with a baseline-penalty term that penalizes visual similarity to the initial state. The resulting reward is continuous, deterministic, and computable in ~ 5 ms per frame, enabling dense per-step feedback during online RL training without task-specific engineering.

We validate Reward-Zero in two stages. First, we develop a completion-sense mini benchmark on ManiSkill simulation tasks, showing that CLIP-direct with baseline penalty achieves 72% forward transition accuracy and perfect jump detection (6/6 episodes), outperforming VLM-caption pipelines (67% best) while being $400\times$ faster. Second, we integrate Reward-Zero as an auxiliary reward into PPO for robotic manipulation and locomotion tasks, demonstrating faster convergence, more stable training dynamics, and higher success rates compared to the PPO baseline with hand-crafted dense rewards alone.

Contribution(s)

1. We propose Reward-Zero, an implicit reward mechanism that uses CLIP vision-language embeddings with a baseline-penalty potential to produce dense completion-sense signals from a natural-language goal description and raw visual observations, without task-specific reward engineering.

Context: Prior language-guided reward methods rely on VLM captioning (~ 2 s/frame) or LLM reward-code synthesis; Reward-Zero operates via direct embedding comparison (~ 5 ms/frame).

2. We introduce a completion-sense mini benchmark that evaluates whether language-grounded reward models assign monotonically increasing potentials across task-completion stages, and use it to show that CLIP-direct with baseline penalty (72% forward accuracy, perfect jump detection) outperforms VLM-caption pipelines (67% best) while being $400\times$ faster.

Context: This benchmark isolates reward-signal fidelity from RL optimization dynamics, which existing evaluations typically conflate.

3. We show that Reward-Zero, integrated as an auxiliary reward into PPO, can accelerate convergence and improve success rates on ManiSkill robotic tasks compared to the PPO baseline with hand-crafted dense rewards alone.

Context: Same PPO hyperparameters and environment configurations as baselines; the only difference is the addition of the Reward-Zero auxiliary signal.

Reward-Zero: Language Embedding Driven Implicit Reward Mechanisms for Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

1 We introduce Reward-Zero, a general-purpose implicit reward mechanism that trans-
2 forms natural-language task descriptions into dense, semantically grounded progress
3 signals for reinforcement learning (RL). Reward-Zero serves as a simple yet sophis-
4 ticated universal reward function that leverages language embeddings for efficient RL
5 training. By comparing the embedding of a task specification with embeddings der-
6 ived from an agent’s interaction experience, Reward-Zero produces a continuous, se-
7 mantically aligned sense-of-completion signal. This reward supplements sparse or de-
8 layed environmental feedback without requiring task-specific engineering. When in-
9 tegrated into standard RL frameworks, it accelerates exploration, stabilizes training,
10 and enhances generalization across diverse tasks. Empirically, agents trained with
11 Reward-Zero converge faster and achieve higher final success rates than conventional
12 methods such as PPO with common reward-shaping baselines, successfully solving
13 tasks that hand-designed rewards could not in some complex tasks. In addition, we
14 develop a mini benchmark for evaluation of completion sense during task execution via
15 language embeddings. These results highlight the promise of language-driven implicit
16 reward functions as a practical path toward more sample-efficient, generalizable, and
17 scalable RL for embodied agents. Code will be released after peer review.

18 1 Introduction

19 Reinforcement learning (RL) has shown remarkable potential in a broad range of domains, from
20 robotic manipulation [Kalashnikov et al. \(2018\)](#); [Zhang et al. \(2024\)](#) and strategic game playing [Sil-
21 ver et al. \(2016\)](#) to autonomous driving [Kendall et al. \(2019\)](#). Its core promise lies in enabling
22 agents to learn complex behaviors directly through interaction with their environments [Sutton &
23 Barto \(2018\)](#), offering a pathway toward adaptive, intelligent systems applicable to real-world prob-
24 lems. However, the success of RL critically depends on the design of effective reward functions that
25 guide the learning process [Yu et al. \(2025\)](#). Engineering these rewards for non-trivial tasks is often
26 a challenging, time-consuming, and error-prone endeavor. Hand-crafted reward signals may capture
27 only partial aspects of a desired behavior, leading to unintended incentives or misaligned learning
28 objectives. In this age of increasingly complex tasks and open-ended environments [Silver & Sut-
29 ton \(2025\)](#), developing more scalable and generalizable reward mechanisms has become a central
30 challenge for advancing modern RL [Goyal et al. \(2019\)](#).

31 Reward design lies at the heart of RL, as it defines the objective that governs agent [Yu et al. \(2025\)](#).
32 Poorly designed rewards can yield brittle or misaligned policies, drastically slowing convergence
33 or encouraging unintended strategies. In many real-world tasks, obtaining accurate, detailed feed-
34 back is extremely difficult, as environments may provide only sparse or delayed signals, and hand-
35 engineered rewards often fail to reflect nuanced notions of progress or partial success. Consequently,
36 researchers have sought more generalizable paradigms for deriving rewards that are grounded in se-
37 mantic understanding rather than hand-tuned [Adeniji et al. \(2023\)](#).

38 Traditional approaches to improving reward signals include techniques such as reward shaping, cur-
39 riculum learning, and auxiliary objectives. Reward shaping methods, for instance, transform sparse
40 environmental rewards into denser signals by adding heuristic-based intermediate rewards [Wiewiora](#)
41 [\(2003\)](#); [Kim et al. \(2025\)](#). Although this can accelerate learning, such methods often require exten-
42 sive domain knowledge and manual tuning. Potential-based shaping provides theoretical guarantees
43 of policy invariance but offers limited flexibility in capturing complex semantic relationships be-
44 tween actions and goals [Wiewiora \(2003\)](#); [Deng et al. \(2025b\)](#). These challenges highlight the limi-
45 tations of purely task-specific designs and motivate the need for a more general, adaptive approach.
46 Ideally, a universal reward mechanism should provide semantically coherent guidance without heavy
47 engineering effort, enabling agents to learn from natural task descriptions and adapt across diverse
48 environments with minimal redesign.

49 Recent advances in language modeling and representation learning offer a compelling direc-
50 tion [Adeniji et al. \(2023\)](#); [Masadome & Harada \(2025\)](#); [Zhang et al. \(2025\)](#): natural language en-
51 capsulates high-level task semantics and contextual intent, suggesting that linguistic structures could
52 serve as a basis for flexible and informative reward mechanisms applicable across domains. How-
53 ever, existing language-guided reward methods often rely on VLM captioning or LLM reward-code
54 synthesis, which can be computationally expensive and may suffer from issues such as goal-echo
55 bias or brittle grounding. Moreover, these methods typically require additional engineering to ensure
56 that the generated rewards are informative and stable during training.

57 To address these challenges, we propose **Reward-Zero**, see Fig. 1, a general-purpose implicit re-
58 ward mechanism that transforms natural-language task descriptions into continuous progress sig-
59 nals. Inspired by the human ability to intuitively gauge task completion from semantic understand-
60 ing, Reward-Zero leverages pre-trained language embeddings to encode both task instructions and
61 the agent’s ongoing experience, computing a semantically aligned "sense of completion" signal that
62 evolves during learning. This implicit reward serves as a universal auxiliary objective that com-
63 plements or substitutes sparse environmental feedback, removing the need for hand-crafted shaping
64 functions. By grounding reward estimation in language semantics, Reward-Zero captures nuanced
65 progress even without explicit task-specific signals. The resulting approach provides a scalable route
66 for improving exploration efficiency, training stability, and generalization across varied domains
67 from simulated manipulation tasks to embodied, real-world agents interacting in open environments.

68 This paper makes three primary contributions. First, we introduce the Reward-Zero, a language
69 embedding-driven implicit reward mechanism. Second, we propose a mini benchmark designed to
70 quantitatively assess an agent’s "completion sense", how well language-grounded feedback tracks
71 progress during task execution. Third, we present comprehensive empirical evaluations demon-
72 strating that Reward-Zero outperforms conventional algorithms such as PPO and standard reward-
73 shaping baselines in terms of convergence speed, stability, and final performance. Through these
74 contributions, our work highlights the promise of language-embedding implicit rewards as a practi-
75 cal step toward more generalizable and sample-efficient RL.

76 **2 Related Work**

77 **2.1 Potential-Based Reward Function**

78 Recent advances in potential-based reward shaping (PBRS) have extended the original theoretical
79 guarantees to richer, trainable shaping functions and intrinsic-motivation signals, improving sample
80 efficiency in sparse-reward tasks while preserving optimal-policy invariance under broader condi-
81 tions [Ng et al. \(1999\)](#); [Wiewiora \(2003\)](#); [Müller & Kudenko \(2025\)](#); [Deng et al. \(2025b\)](#). Work inte-
82 grating PBRS with intrinsic curiosity and count-based exploration has enabled faster exploration in
83 procedurally complex environments [Pathak et al. \(2017\)](#); [Bellemare et al. \(2016\)](#); [Xu et al. \(2025\)](#).
84 At the same time, practical limitations remain: many modern intrinsic rewards are nonstationary
85 or depend on learned representations, which can violate PBRS assumptions and induce policy bias;
86 computational overhead and sensitivity to potential design also hinder transferability across domains

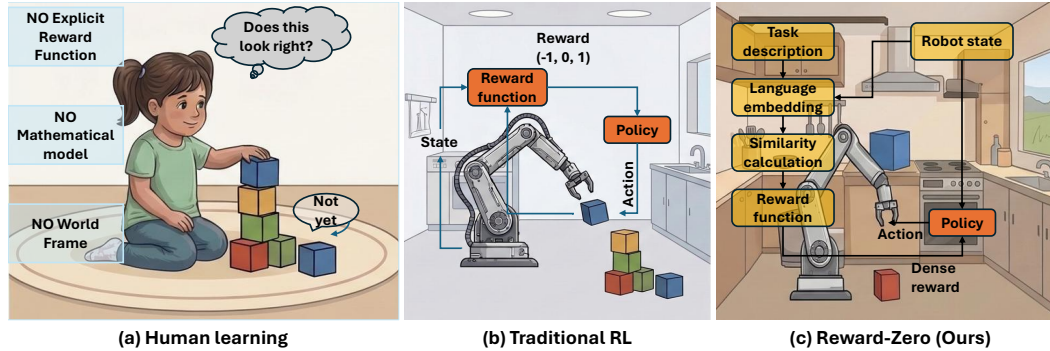


Figure 1: Conceptual comparison of human learning, traditional RL, and the proposed Reward-Zero. The left panel illustrates how Human Learning is intuitive and implicit, driven by visual matching and a generalized "Sense of Completion," without relying on explicit rewards, mathematical environment models, or exact object world coordinates. The center panel depicts Traditional RL training a simple robotic policy, a rigid and explicit approach that typically requires a hand-crafted reward function, a precise environment model, heavy observation and exact object coordinates to function. The right panel introduces Reward-Zero, our flexible and language-driven method that represents a sophisticated universal reward function. This mechanism uses a general-purpose language embedding-driven implicit reward mechanism to generate a continuous sense-of-completion signal by comparing task and experience embeddings. As shown, Reward-Zero aims to eliminate hand-crafted rewards by relying only on raw language embedding, enhancing generalization across diverse tasks. **Zero** here signifies the absence of hand-crafted rewards, without explicit reward engineering. This is the **Zero** step toward more general, adaptable, and scalable RL that can learn from natural language descriptions and raw observations, much like humans do.

87 [Sutton & Barto \(2018\)](#); [Farooq et al. \(2025\)](#); [Houthoofd et al. \(2016\)](#); [Müller & Kudenko \(2025\)](#). Recent methods therefore focus on converting learned intrinsic signals into potential-based forms or
 88 constraining their dependence to preserve optimality, but trade-offs between theoretical guarantees
 89 and empirical performance persist [Kim et al. \(2025\)](#); [Deng et al. \(2025b\)](#).
 90

91 2.2 Language-Guided Reward for RL

92 Language-guided reward learning uses natural-language specifications and large language models
 93 to produce dense, interpretable reward signals that accelerate learning in sparse, compositional, or
 94 multi-step tasks. Recent work demonstrates automatic synthesis of shaped rewards from textual
 95 goals, enabling data-efficient policy learning and iterative human refinement for robotics and em-
 96 bodied benchmarks [Xie et al. \(2023\)](#); [Zhang et al. \(2025\)](#). Integrations that fine-tune or prompt
 97 LLMs to generate reward functions or reward classifiers show improved generalization across task
 98 variants and reduce the need for new demonstrations [Deng et al. \(2025a\)](#); [Adeniji et al. \(2023\)](#).
 99 Other lines pretrain joint policy–reward models on language-annotated trajectories to transfer to
 100 unseen instructions [Masadome & Harada \(2025\)](#). Remaining challenges include brittle grounding
 101 when environment observations misalign with textual abstractions, latent biases in LLM-derived re-
 102 wards, and the need for verification to prevent reward hacking; recent proposals therefore emphasize
 103 human-in-the-loop refinement, robust grounding mechanisms, and formalizing safety guarantees for
 104 language-derived rewards [Zhang et al. \(2025\)](#); [Ji et al. \(2026\)](#); [Adeniji et al. \(2023\)](#).

105 3 Method

106 In this section, we present Reward-Zero, a vision-language model-based reward computation ap-
 107 proach that provides dense, semantically grounded rewards for robotic manipulation tasks. Our

108 method consists of three key components: (1) language embedding-based potential estimation, (2)
109 progress-aware activation, and (3) completion-sense reward shaping.

110 3.1 Language Embedding-Based Potential Estimation

111 Traditional reward shaping in robotics often relies on hand-crafted distance metrics or task-specific
112 state features, which limits generalization across diverse manipulation tasks. We propose a fun-
113 damentally different approach that leverages the semantic understanding capabilities of vision-
114 language models (VLMs) to compute task-relevant potentials directly from natural language de-
115 scriptions.

116 Our core insight is that the semantic similarity between a scene description and a goal description
117 naturally captures task progress without requiring explicit geometric or kinematic computations.
118 Specifically, we define the potential function $\Phi(s)$ as the cosine similarity between the encoded
119 current state caption and the encoded goal description:

$$\Phi(s) = \cos(\mathbf{e}_{\text{state}}, \mathbf{e}_{\text{goal}}) = \frac{\mathbf{e}_{\text{state}} \cdot \mathbf{e}_{\text{goal}}}{\|\mathbf{e}_{\text{state}}\| \|\mathbf{e}_{\text{goal}}\|} \quad (1)$$

120 where $\mathbf{e}_{\text{state}}$ and \mathbf{e}_{goal} are the text embeddings of the current scene caption and goal description,
121 respectively. This formulation yields a bounded potential $\Phi(s) \in [-1, 1]$, where higher values
122 indicate greater semantic alignment with the goal.

123 To enhance the expressiveness and discriminability of our potential function, we employ an enrich-
124 ment procedure for both state captions and goal descriptions. Rather than using raw, terse captions
125 (e.g., “robot arm near cup”), we leverage large language models to generate detailed, contextually-
126 rich descriptions that capture nuanced aspects of the scene. For state descriptions, the VLM cap-
127 tioner is prompted to produce comprehensive scene descriptions that include object positions, spa-
128 tial relationships, gripper states, and ongoing actions. Similarly, goal descriptions are enriched to
129 include expected final configurations, success criteria, and relevant contextual details.

130 This enrichment process is crucial for two reasons. First, richer text provides more distinctive em-
131 beddings in the semantic space, enabling finer-grained progress discrimination. Second, detailed
132 descriptions help bridge the semantic gap between visual observations and abstract task specifica-
133 tions, allowing the embedding space to capture task-relevant features more effectively. The general-
134 ity of this approach stems from its reliance solely on pretrained language models and text encoders,
135 requiring no task-specific engineering or domain knowledge.

136 3.2 Progress-Aware Activation

137 While the potential function provides a continuous measure of goal proximity, raw potential differ-
138 ences may not adequately incentivize the agent during critical final stages of task completion. We
139 introduce a progress-aware activation mechanism that dynamically amplifies rewards as the agent
140 approaches task completion.

141 The activation function is based on a sigmoid transformation centered at a completion threshold τ :

$$\sigma_{\text{act}}(\Phi) = \frac{1}{1 + \exp(-k \cdot (\Phi - \tau))} \quad (2)$$

142 where k controls the steepness of the activation transition and τ represents the potential value at
143 which the agent is considered to be near completion. This sigmoid activation provides several desir-
144 able properties: (1) it remains near zero when the agent is far from the goal, avoiding interference
145 with exploration; (2) it smoothly transitions to high values as the agent approaches completion,
146 providing increasingly strong guidance; and (3) it avoids discontinuous reward jumps that could
147 destabilize learning.

148 To further encourage sustained progress during the critical completion phase, we incorporate a
 149 progress multiplier that rewards continued improvement:

$$\Delta\Phi = \max(0, \Phi_t - \Phi_{t-1}) \quad (3)$$

150 The progress term ensures that the agent receives additional reward for making forward progress
 151 even when already close to the goal. This addresses a common challenge in potential-based shaping
 152 where the reward signal diminishes as the agent approaches the goal state, potentially leading to
 153 sluggish final movements. By multiplying the activation by $(1 + \Delta\Phi)$, we create a reward landscape
 154 that actively pulls the agent toward completion rather than merely indicating proximity.

155 3.3 Completion-Sense Reward Formulation

156 Combining the language-based potential with progress-aware activation, we formulate the complete
 157 Reward-Zero function as:

$$R_{\text{completion}} = r_{\text{base}} + \beta \cdot \sigma_{\text{act}}(\Phi) \cdot (1 + \Delta\Phi) \quad (4)$$

158 where r_{base} represents the base potential reward (which can be Φ itself or the potential difference
 159 $\Phi_t - \Phi_{t-1}$), β is the completion bonus weight, $\sigma_{\text{act}}(\Phi)$ is the sigmoid activation function, and $\Delta\Phi$
 160 captures the instantaneous progress.

161 This formulation achieves several design objectives that make it effective for general robotic manip-
 162 ulation:

163 **Smoothness and Stability:** Unlike sparse rewards or threshold-based completion bonuses, our re-
 164 ward function is continuous and differentiable everywhere. This property is essential for stable
 165 policy gradient estimation and avoids the optimization difficulties associated with discontinuous
 166 reward landscapes.

167 **Automatic Curriculum:** The reward naturally implements a form of curriculum learning. Early in
 168 training, when the agent is far from the goal, the completion bonus remains inactive, allowing the
 169 base potential to guide exploration. As the agent learns to approach the goal, the completion bonus
 170 gradually activates, providing stronger incentives for precise final positioning.

171 **Task Generalization:** By grounding rewards in language semantics rather than task-specific met-
 172 rics, our approach generalizes across diverse manipulation tasks without modification. The same
 173 reward function can evaluate progress toward “pick up the red cube,” “close the drawer,” or “stack
 174 blocks in order” simply by changing the goal text description.

175 **Dense Feedback:** Every timestep receives meaningful reward signal derived from semantic scene
 176 understanding, addressing the sparse reward challenge that plagues many manipulation tasks. This
 177 density accelerates learning and provides consistent gradient information throughout training.

178 The hyperparameters τ , k , and β can be tuned based on task characteristics, though we find that
 179 default values of $\tau = 0.7$, $k = 10$, and $\beta = 0.5$ work well across a variety of manipulation
 180 scenarios. The threshold τ should be set below the expected completion potential to ensure the
 181 bonus activates during the approach phase rather than only at completion.

182 4 Experiments and Results

183 In this section, we evaluate the effectiveness of our proposed Reward-Zero across a range of robotic
 184 manipulation tasks. We compare our method against standard RL algorithms and common reward-
 185 shaping baselines to demonstrate its advantages in terms of sample efficiency, convergence speed,
 186 and final performance (success rate). First, we develop a mini benchmark for evaluation of com-
 187 pletion sense during task execution via language embeddings, see Section 4.1. Then, we evaluate

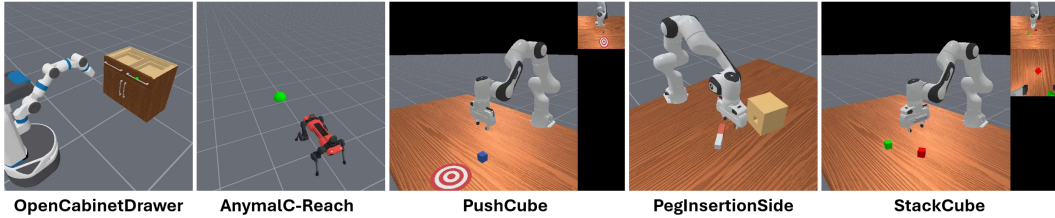


Figure 2: Example tasks and keyframes from the completion-sense mini benchmark. Each episode contains 2–4 annotated keyframes at known completion percentages (0%, 33%, 50%, 66%, 100%) extracted from successful ManiSkill Gu et al. (2023) rollouts. The benchmark includes tasks with varying visual complexity, from large state changes (e.g., OpenCabinetDrawer) to fine-grained manipulations (e.g., PegInsertionSide).

188 our method on a suite of robotic manipulation tasks that require varying levels of complexity and
 189 dynamic interaction, see Section 4.2.

190 4.1 Mini Benchmark for Completion-Sense Evaluation

191 Before deploying Reward-Zero as an auxiliary reward signal in online RL training, we must verify
 192 that the underlying potential function $\Phi(s)$ reliably tracks task-completion progress from visual ob-
 193 servations alone. To this end, we develop a lightweight offline benchmark that evaluates *completion-*
 194 *sense discrimination*: the ability of a reward model to assign monotonically increasing potential val-
 195 ues to frames sampled at increasing stages of task completion, given only a natural-language goal
 196 description.

197 **Benchmark Design** We construct evaluation episodes from ManiSkill Gu et al. (2023) simulation
 198 trajectories across five robotic tasks spanning locomotion, tabletop manipulation, and articulated-
 199 object interaction (Fig. 2 and Tab. 1). For each task, we extract keyframes at four known com-
 200 pletion percentages (0%, 33%, 66%, 100%) from successful rollouts, yielding 4 annotated frames
 201 per episode. Each task is paired with an end-state goal description (e.g., “The cabinet drawer is
 202 fully open”) rather than an action command, as we found that action-phrased goals (e.g., “Open the
 203 drawer by pulling the handle”) cause VLM-based models to echo goal language even at 0% com-
 204 pletion, inflating similarity scores. In total, the benchmark contains 6 evaluation episodes across 5
 205 environments with 24 annotated keyframes, providing 18 consecutive forward transitions for assess-
 206 ment.

Table 1: Tasks in the completion-sense mini benchmark. Each episode contains keyframes at known completion percentages extracted from successful ManiSkill rollouts. Two episodes are included for OpenCabinetDrawer to test consistency across different initial configurations.

Task	Frames	Completions (%)	Goal Description
OpenCabinetDrawer ($\times 2$)	4	0, 33, 66, 100	The cabinet drawer is fully open
AnymalC-Reach	4	0, 33, 66, 100	The quadruped robot is at the target position
PushCube	4	0, 33, 66, 100	The cube is at the target position
PegInsertionSide	4	0, 33, 66, 100	The peg is fully inserted into the hole
StackCube	4	0, 33, 66, 100	The red cube is stacked on the green cube

207 These tasks were selected to span a range of visual complexities: OpenCabinetDrawer and
 208 AnymalC-Reach involve large, visually distinctive state changes, while PegInsertionSide and Stack-
 209 Cube involve fine-grained manipulations where objects are small relative to the scene.

210 **Evaluation Metrics.** We evaluate reward models along four complementary dimensions:

- 211 • **Forward Transition Accuracy (FTA):** For each consecutive pair of frames (s_t, s_{t+1}) where
 212 completion increases, we check whether the computed reward $R(s_t, s_{t+1}) > \epsilon$ (with threshold
 213 $\epsilon = 0.001$). This directly measures whether the reward signal provides a positive learning
 214 gradient in the direction of task progress. Reported out of 18 total transitions.
- 215 • **Monotonicity:** The fraction of consecutive potential pairs $(\Phi(s_t), \Phi(s_{t+1}))$ that strictly in-
 216 crease. A score of 1.0 indicates perfectly monotonic potential tracking across the episode.
- 217 • **Spearman Correlation (ρ):** Rank correlation between ground-truth completion percentages and
 218 computed potential values, capturing ordinal alignment without assuming linearity.
- 219 • **Jump Detection (J+):** Whether a single 0%→100% transition produces a clearly positive reward
 220 ($R > \epsilon$), testing the model’s sensitivity to large state changes. Reported out of 6 episodes.

221 Forward transition accuracy and jump detection are the most directly relevant metrics for RL train-
 222 ing, where the agent always progresses forward from an initial state.

223 **Methods Compared.** We evaluate two families of approaches for instantiating the potential func-
 224 tion $\Phi(s)$:

225 (1) *VLM-caption pipeline.* An image is passed through a vision-language model (VLM) which
 226 generates a natural-language scene description. This caption is encoded using a sentence embedding
 227 model (MiniLM-L6 Wang et al. (2020)), and the cosine similarity between the caption embedding
 228 and the goal-text embedding serves as the potential. We test Qwen2.5-VL-3B Bai et al. (2025) with
 229 three prompting strategies:

- 230 • **Progress-description:** The prompt provides the goal text and asks the VLM to describe both the
 231 current scene state and progress toward the goal. This is the most informative prompt but risks
 232 goal-echo bias.
- 233 • **Observe-only:** The prompt asks the VLM to describe the scene without revealing the goal,
 234 avoiding echo bias at the cost of less goal-directed captions.
- 235 • **Evidence-gating:** The prompt instructs the VLM to enumerate only actions that have been
 236 visibly completed, with a fallback phrase (“No actions completed”) if uncertain.

237 (2) *CLIP-direct.* The image is encoded directly by CLIP’s Radford et al. (2021) vision encoder (ViT-
 238 B/32) and compared against the CLIP text encoding of the goal description. No VLM captioning or
 239 intermediate text representation is needed. We define the potential as:

$$\Phi(s) = \alpha \cdot \text{sim}(f_I(s), f_T(g)) - (1 - \alpha) \cdot \text{sim}(f_I(s), f_I(s_0)) \quad (5)$$

240 where f_I and f_T are CLIP’s image and text encoders, g is the goal text, and s_0 is the initial obser-
 241 vation (episode baseline). The first term measures goal proximity; the second term penalizes visual
 242 similarity to the initial state, encouraging departure from the start configuration. We use $\alpha = 0.7$,
 243 balancing goal affinity with departure from the initial state.

244 **Results.** Table 2 presents the main comparison across all methods. CLIP-direct with baseline
 245 penalty ($\alpha = 0.7$) achieves the highest forward transition accuracy (13/18, 72%) and perfect jump
 246 detection (6/6), outperforming all VLM-based approaches.

247 We showcase the per-task keyframes and discussion in Appendix A.

248 **Key Findings.** Several findings emerge from the aggregate comparison (Table 2) and per-task
 249 analysis (Figures 9–3):

250 **CLIP-direct outperforms VLM-caption pipelines on forward discrimination.** The best VLM
 251 configurations (progress-description and observe-only) each achieve 67% forward accuracy (12/18),
 252 while CLIP-direct with $\alpha = 0.7$ reaches 72% (13/18). More importantly, CLIP achieves per-
 253 fect jump detection (6/6) and is 400× faster. The VLM-caption pipeline introduces noise at two
 254 stages: first, VLM-generated captions frequently hallucinate task progress (e.g., describing a peg

Table 2: Mini benchmark results comparing VLM-caption and CLIP-direct approaches on 18 forward transitions across 6 episodes. **FTA**: forward transition accuracy (out of 18). **J+**: jump detection, 0%→100% (out of 6). **Mono**: fraction of episodes with fully monotonic potentials (out of 6). Latency is per-frame inference time on a single A30 GPU. Best results in **bold**. CLIP results are deterministic.

Method	FTA (\uparrow)	J+ (\uparrow)	Mono (\uparrow)	Latency
<i>VLM-caption pipeline (Qwen2.5-VL-3B + MiniLM)</i>				
+ progress-description	12/18	5/6	2/6	~2 s
+ observe-only	12/18	6/6	1/6	~2 s
+ evidence-gating	0/18	0/6	0/6	~2 s
<i>CLIP-direct (ViT-B/32), $\alpha = 0.7$</i>				
with baseline penalty	13/18	6/6	2/6	~5 ms

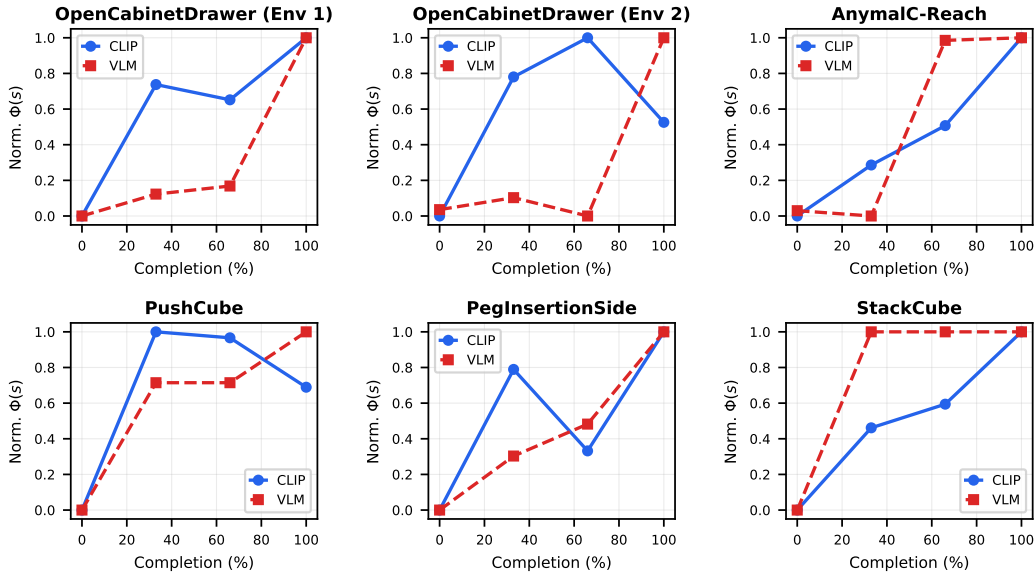


Figure 3: CLIP potential $\Phi(s)$ vs. task completion (%) for each benchmark task (CLIP-direct, $\alpha = 0.7$). OpenCabinetDrawer shows two episodes (solid/dashed) to illustrate consistency across initial configurations. All tasks now use four keyframes at 0%, 33%, 66%, and 100% completion. Tasks with large visual changes show strong monotonic trends; fine-manipulation tasks exhibit smaller potential ranges.

255 as “partially inserted” when it lies on the table at 0% completion); second, the additional sentence-
 256 embedding step (MiniLM) maps semantically distinct captions to similar embedding vectors, further
 257 diluting the discrimination signal.

258 **The baseline penalty term is theoretically motivated and practically beneficial.** The term
 259 $-(1 - \alpha) \cdot \text{sim}(f_I(s), f_I(s_0))$ penalizes states that remain visually similar to the initial frame, am-
 260 plying the reward gradient in the forward direction. This introduces an expected asymmetry where
 261 backward transitions receive less reward, which is acceptable for RL training since episodes always
 262 progress forward from the initial state.

263 **Completion-sense quality correlates with visual saliency.** Tasks where the goal state is visually
 264 distinctive from the initial state (large object displacement, prominent structural changes) produce
 265 strong, monotonic potential trajectories. Fine-grained manipulation tasks, where relevant objects
 266 are small or partially occluded, yield noisier signals—motivating the use of higher-resolution vision
 267 encoders (e.g., ViT-B/16) for such domains.

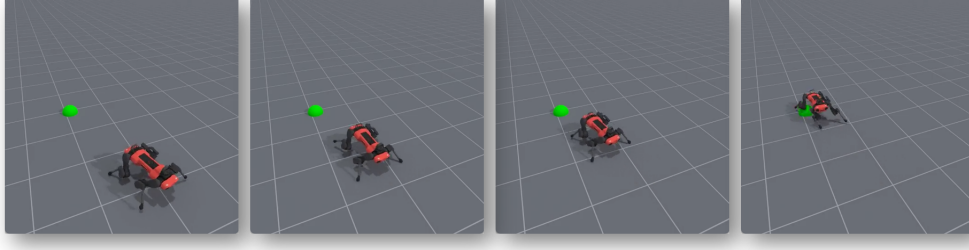


Figure 4: The AnymalC-Reach task involves a quadruped robot learning to navigate to a target location. The task requires the agent to understand spatial relationships and adapt its locomotion strategy accordingly. This environment serves as a challenging testbed for evaluating the effectiveness of our Reward-Zero in guiding learning through language-driven implicit rewards.

268 **VLM prompting strategy has a dramatic effect.** The evidence-gating prompt achieves 0/18 accu-
 269 racy because it provides a safe fallback phrase (“No actions completed”), which the VLM defaults
 270 to for every frame regardless of actual completion. Progress-description and observe-only prompts
 271 achieve identical FTA (12/18), though observe-only achieves perfect jump detection (6/6) by de-
 272 scribing scene state more faithfully without goal-echoing bias.

273 **CLIP is deterministic and orders of magnitude faster.** CLIP processes a frame in ~ 5 ms on GPU,
 274 while VLM captioning requires ~ 2 s per frame—a $400\times$ speed difference. This enables dense per-
 275 step reward computation during RL training, whereas VLM-based rewards must be sparsely sampled
 276 (e.g., every 25–50 steps), further disadvantaging VLM approaches for online use.

277 Based on these results, we adopt the CLIP-direct formulation (Eq. 5) with $\alpha = 0.7$ as the default
 278 Reward-Zero potential function for all subsequent RL experiments.

279 4.2 Embodied Tasks

280 To further deeply evaluate the effectiveness of our proposed Reward-Zero, we conduct experiments
 281 on a suite of robotic tasks that require varying levels of complexity and dynamic interactions. These
 282 tasks are designed to test the generalization capabilities of our method across different scenarios and
 283 to compare its performance against standard RL algorithms and common reward-shaping baselines.
 284 One example task is the quadruped robot “AnymalC-Reach” task shown in Fig. 4, where the robot
 285 must learn to navigate to a target location.

286 We compare the performance of Reward-Zero against traditional RL methods (PPO) and reward-
 287 shaping baselines that rely on hand-crafted distance metrics or task-specific features. Our evaluation
 288 metrics include convergence speed, sample efficiency, and final success rate. The results in Fig. 5
 289 demonstrate that Reward-Zero significantly accelerates learning and improves generalization across
 290 tasks compared to traditional approaches, highlighting the advantages of our language-driven im-
 291 plicit reward mechanism in guiding agents toward task completion in a more semantically meaning-
 ful way.

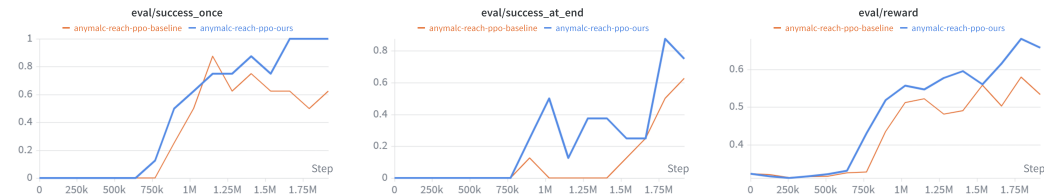


Figure 5: Performance comparison of the ANYmal-C Reach-PPO baseline versus our approach. Solid lines represent the mean values over 2M training steps. Our method significantly outperforms the baseline in both task success rates and cumulative reward.

292 Moreover, we analyze the learning curves in Fig. 6 demonstrating that ours yields substantially
 293 more stable and reliable optimization dynamics compared to the PPO baseline. The value loss of
 294 the baseline exhibits pronounced oscillations, particularly in later stages, indicating unstable critic
 295 fitting, whereas Reward-Zero maintains a consistently smooth trajectory, reflecting a more accurate
 296 and stable value function. Similarly, the policy-related metrics, including policy loss, KL diver-
 297 gence, and clipping fraction, show that Reward-Zero performs updates with markedly fewer spikes
 298 and reduced variance, suggesting more controlled policy improvement and a lower risk of cata-
 299 strophic updates. In terms of learning efficiency, Reward-Zero achieves faster entropy reduction and
 300 maintains higher, more stable explained variance, implying quicker convergence toward effective
 301 policies and more reliable advantage estimation. Collectively, these results indicate that Reward-
 302 Zero not only stabilizes PPO training but also improves sample efficiency and robustness, ultimately
 303 enabling more consistent and higher-quality policy learning.

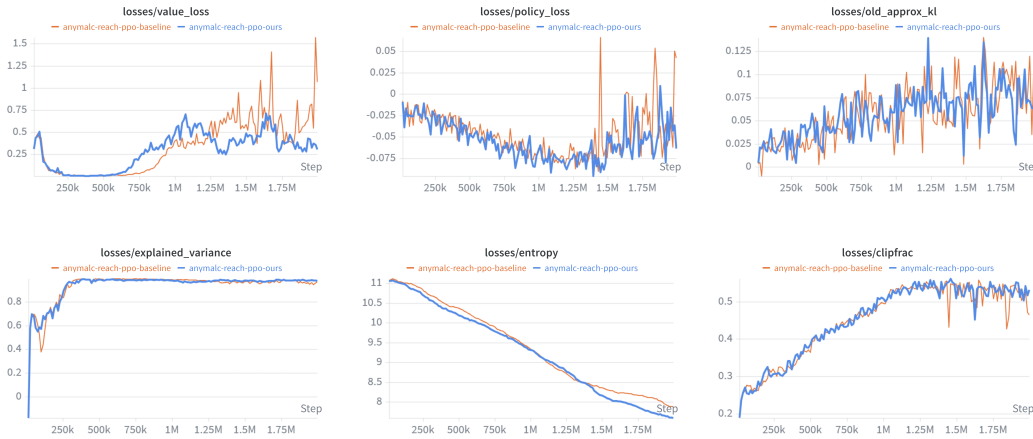


Figure 6: Training dynamics comparison between the PPO baseline (orange) and Ours (blue). The Reward-Zero variant exhibits substantially smoother value loss, more stable policy updates, and more consistent KL and clipping behavior, indicating improved critic reliability and better-controlled policy optimization. Reward-Zero also shows faster entropy reduction and more stable explained variance, reflecting higher learning efficiency and more reliable advantage estimation throughout training.

304 **4.3 Ablation Study**

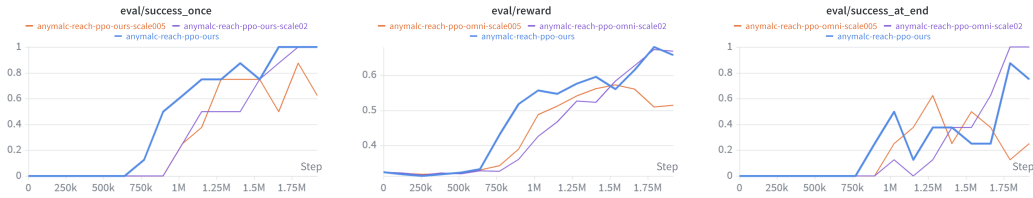


Figure 7: Ablation study on the scale parameter for Reward-Zero policy learning.

305 We study the effect of the scale parameter β in the Reward-Zero formulation (Eq. 5) on policy
 306 learning performance. In Fig. 7, we compare three configurations: $\beta = 0.05$ (low completion bonus,
 307 orange), $\beta = 0.1$ (default as ours, blue), and $\beta = 0.2$ (high completion bonus, purple). The results
 308 across success-once, success-at-end, and the reward. While reducing or increasing the
 309 scale parameter alters learning dynamics, the blue curve (our default) consistently achieves the most
 310 reliable and highest overall performance. It reaches perfect success-once early, maintains strong
 311 success-at-end behavior, and yields the highest evaluation reward, indicating that the chosen

312 scale value provides the best balance between exploration magnitude and policy update stability.
 313 These results highlight that proper scale calibration is crucial for maximizing the effectiveness of
 completion-sense rewards and ensuring robust task completion.

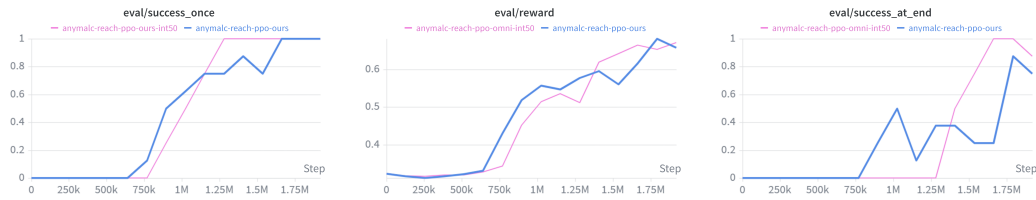


Figure 8: Ablation on Reward-Zero Invocation Frequency.

314

315 We further examine how frequently Reward-Zero should be invoked during training by comparing
 316 our default configuration ($interval = 25$) with a less frequent variant ($interval = 50$). As shown
 317 in Fig. 8, the results reveal a clear trade-off between reward-signal density and policy stability,
 318 and they highlight why the default setting provides the most balanced and reliable performance.
 319 The evaluation reward curve further clarifies this distinction. Although both variants improve over
 320 time, the default interval ultimately achieves slightly higher reward, reflecting more stable policy
 321 refinement and better alignment between intermediate progress signals and final task success.

322 5 Conclusions and Future Work

323 In this paper, we introduced Reward-Zero, a novel reward mechanism for computing dense,
 324 semantically-grounded rewards. By leveraging the semantic similarity between scene descriptions
 325 and goal specifications, our method provides a powerful potential function that guides learning with-
 326 out requiring task-specific engineering. We demonstrated the effectiveness of Reward-Zero through
 327 a mini benchmark evaluating completion-sense discrimination and through experiments on complex
 328 robotic tasks, showing improvements in sample efficiency and convergence speed compared to tradi-
 329 tional RL approaches. These results underscore the promise of language-driven reward shaping
 330 for sparse-reward settings, enabling more generalizable and interpretable reward functions. Future
 331 directions include fully language-embedding-based reward models and success criteria, and deploy-
 332 ment in real-world robotic systems.

333 References

- 334 Ademi Adeniji, Amber Xie, Carmelo Sferrazza, Younggyo Seo, Stephen James, and Pieter
 335 Abbeel. Language reward modulation for pretraining reinforcement learning. *arXiv preprint*
 336 *arXiv:2308.12270*, 2023.
- 337 Shuai Bai, Yiheng Xu, Peng Wang, Hang Zhang, Pengfei Wang, Shijie Wang, Junyang Lin,
 338 Tianbao Xie, Yuanzhi Zhu, Zhibo Yang, et al. Qwen2.5-vl technical report. *arXiv preprint*
 339 *arXiv:2502.13923*, 2025.
- 340 Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi
 341 Munos. Unifying count-based exploration and intrinsic motivation. *Advances in Neural Informa-*
 342 *tion Processing Systems*, 29, 2016.
- 343 Yongxin Deng, Xihe Qiu, Jue Chen, and Xiaoyu Tan. Reward guidance for reinforcement learning
 344 tasks based on large language models: The limgt framework. *Knowledge-Based Systems*, 322:
 345 113689, 2025a.
- 346 Zelin Deng, Yunlong Dong, and Xing Liu. Reward shaping in reinforcement learning for robotic
 347 hand manipulation. *Neurocomputing*, 638:130204, 2025b.

- 348 Sehar Shahzad Farooq, Hameedur Rahman, Samiya Abdul Wahid, Muhammad Alyan Ansari, Saira
349 Abdul Wahid, and Hosu Lee. Curiosity-driven exploration in reinforcement learning: An adaptive
350 self-supervised learning approach for playing action games. *Computers*, 14(10):434, 2025.
- 351 Prasoon Goyal, Scott Niekum, and Raymond J Mooney. Using natural language for reward shaping
352 in reinforcement learning. In *Proceedings of the 28th International Joint Conference on Artificial
353 Intelligence*, pp. 2385–2391, 2019.
- 354 Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhanjia Ling, Xuan Liu, Tongzhou Mu, Yihe Tang, Stone Tao,
355 Xinyue Wei, Yuzhe Yao, et al. Maniskill2: A unified benchmark for generalizable manipulation
356 skills. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL
357 https://openreview.net/forum?id=S7cl9AtY_79.
- 358 Rein Houthoofd, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime:
359 Variational information maximizing exploration. In *Advances in Neural Information Processing
360 Systems*, 2016.
- 361 Miaomiao Ji, Yanqiu Wu, Zhibin Wu, Shoujin Wang, Jian Yang, Mark Dras, and Usman Naseem.
362 A survey of progress in llm alignment from the perspective of reward design. *IEEE Transactions
363 on Artificial Intelligence*, 2026.
- 364 Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre
365 Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforce-
366 ment learning for vision-based robotic manipulation. In *Conference on robot learning*, pp. 651–
367 673. PMLR, 2018.
- 368 Alex Kendall, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John-Mark Allen,
369 Vinh-Dieu Lam, Alex Bewley, and Amar Shah. Learning to drive in a day. In *2019 International
370 Conference on Robotics and Automation (ICRA)*, pp. 8248–8254, 2019. DOI: 10.1109/ICRA.
371 2019.8793742.
- 372 Dohyeong Kim, Hyeokjin Kwon, Junseok Kim, Gunmin Lee, and Songhwai Oh. Stage-wise reward
373 shaping for acrobatic robots: A constrained multi-objective reinforcement learning approach. In
374 *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 10268–10274,
375 2025. DOI: 10.1109/ICRA55743.2025.11128552.
- 376 Shinya Masadome and Taku Harada. Reward design using large language models for natural lan-
377 guage explanation of reinforcement learning agent actions. *IEEJ Transactions on Electrical and
378 Electronic Engineering*, 20(8):1203–1211, 2025.
- 379 Henrik Müller and Daniel Kudenko. Improving the effectiveness of potential-based reward shaping
380 in reinforcement learning. *arXiv preprint arXiv:2502.01307*, 2025.
- 381 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:
382 Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287. Citeseer, 1999.
- 383 Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration
384 by self-supervised prediction. In *ICML Deep RL Workshop*, 2017.
- 385 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
386 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
387 models from natural language supervision. In *International Conference on Machine Learning
388 (ICML)*, pp. 8748–8763. PMLR, 2021.
- 389 David Silver and Richard S Sutton. Welcome to the era of experience. *Google AI*, 1, 2025.
- 390 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
391 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
392 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

- 393 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd
394 edition, 2018.
- 395 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-
396 attention distillation for task-agnostic compression of pre-trained transformers. *Advances in neu-
397 ral information processing systems*, 33:5776–5788, 2020.
- 398 Eric Wiewiora. Potential-based shaping and q-value initialization are equivalent. *Journal of Artificial
399 Intelligence Research*, 19:205–208, 2003.
- 400 Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang, and
401 Tao Yu. Text2reward: Reward shaping with language models for reinforcement learning. *arXiv
402 preprint arXiv:2309.11489*, 2023.
- 403 Feng Xu, Guangyao Zhai, Xin Kong, Tingzhong Fu, Daniel FN Gordon, Xueli An, and Benjamin
404 Busam. Stare-vla: Progressive stage-aware reinforcement for fine-tuning vision-language-action
405 models. *arXiv preprint arXiv:2512.05107*, 2025.
- 406 Rui Yu, Shenghua Wan, Yucen Wang, Chen-Xiao Gao, Le Gan, Zongzhang Zhang, and De-Chuan
407 Zhan. Reward models in deep reinforcement learning: a survey. In *Proceedings of the Thirty-
408 Fourth International Joint Conference on Artificial Intelligence*, pp. 10807–10816, 2025.
- 409 Heng Zhang, Gokhan Solak, Gustavo J. G. Lahr, and Arash Ajoudani. Srl-vic: A variable stiffness-
410 based safe reinforcement learning for contact-rich robotic tasks. *IEEE Robotics and Automation
411 Letters*, 9(6):5631–5638, 2024. DOI: 10.1109/LRA.2024.3396368.
- 412 Jiahui Zhang, Yusen Luo, Abrar Anwar, Sumedh Anand Sontakke, Joseph J Lim, Jesse Thomason,
413 Erdem Biyik, and Jesse Zhang. Rewind: Language-guided rewards teach robot policies without
414 new demonstrations. *arXiv preprint arXiv:2505.10911*, 2025.

415
416
417

Supplementary Materials

The following content was not necessarily subject to peer review.

418 A Mini benchmark for evaluation of completion sense during task execution

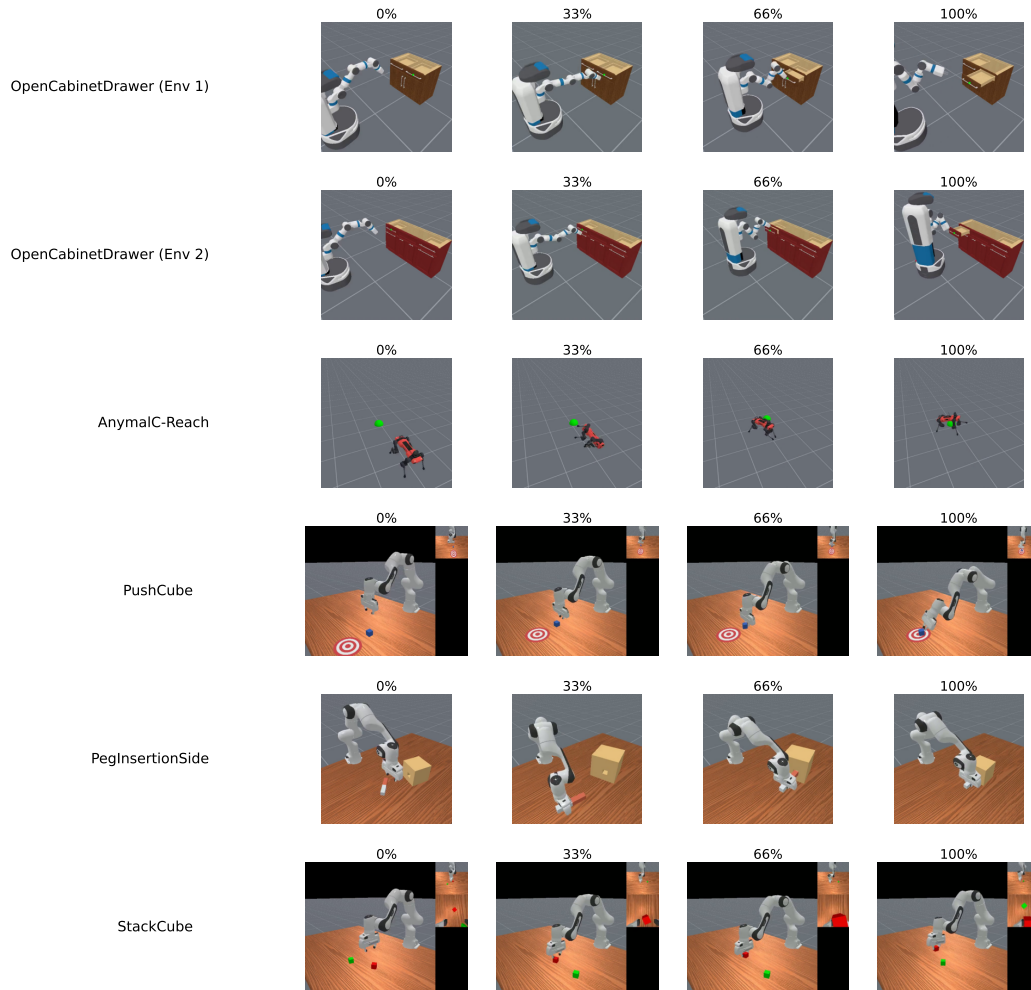


Figure 9: Keyframes from the six benchmark episodes (five task types) at four completion stages (0%, 33%, 66%, 100%). Each row shows one episode; columns correspond to the annotated completion percentages. Tasks range from visually salient changes (top: drawer opening, robot locomotion) to fine-grained manipulation where task-relevant objects are small relative to the scene (bottom: peg insertion, cube stacking).

419 Figure 9 shows representative keyframes for all five tasks at increasing completion stages, and Fig-
 420 ure 3 plots the corresponding CLIP potential $\Phi(s)$ curves using our best configuration ($\alpha = 0.7$).
 421 The tasks span a range of visual complexities: OpenCabinetDrawer and AnymalC-Reach involve
 422 large, structurally distinctive state changes and produce reliably monotonic potentials, while fine-
 423 manipulation tasks (PegInsertionSide, StackCube) involve small objects that occupy only a few
 424 pixels at CLIP’s 224×224 input resolution, yielding smaller dynamic ranges. Notably, PegInser-
 425 tionSide’s peg (~ 2 cm physical, ~ 7 px in CLIP input) fits within a single ViT-B/32 patch, limiting
 426 CLIP’s ability to track fine-grained insertion progress.