AUTOMATED REWARDS VIA LLM-GENERATED PROGRESS FUNCTIONS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have the potential to automate reward engineering for Reinforcement Learning (RL) by leveraging their broad domain knowledge across various tasks. However, they often need many iterations of trial-and-error to generate effective reward functions. This process is costly because evaluating every sampled reward function requires completing the full policy optimization process for each function. In this paper, we introduce an LLM-driven reward generation framework that is able to produce state-of-the-art policies on the challenging Bi-DexHands benchmark with 20× fewer reward function samples than the prior state-of-the-art work. Our key insight is that we reduce the problem of generating task-specific rewards to the problem of coarsely estimating task progress. Our two-step solution leverages the task domain knowledge and the code synthesis abilities of LLMs to author *progress functions* that estimate task progress from a given state. Then, we use this notion of progress to discretize states, and generate count-based intrinsic rewards using the low-dimensional state space. We show that the combination of LLM-generated progress functions and count-based intrinsic rewards is essential for our performance gains, while alternatives such as generic hash-based counts or using progress directly as a reward function fall short.

028 1 INTRODUCTION 029

Automated reward engineering aims to reduce the human effort required when using Reinforcement
 Learning (RL) for sparse-reward tasks. Multiple recent efforts towards automated reward engineering
 have focused on leveraging Large Language Models (LLMs) to provide reward signals—either
 employing the LLM output directly as the reward (Kwon et al., 2023), or using the LLM to generate
 code for a dense reward function (Yu et al., 2023; Ma et al., 2023).

Traditionally, constructing an effective dense reward function is an intricate process of identifying key task elements and carefully weighing different reward terms (Sutton & Barto, 2018; Booth et al., 2023). Prior works on reward function generation attempt to use LLMs both for their domain knowledge and to optimize quantitative aspects of reward engineering (weighting, rescaling) (Yu et al., 2023; Ma et al., 2023)–they can require many training runs as they search through many reward functions in order to find a candidate that is effective for training (Ma et al., 2023).

041 Our core insight is that we can reduce the problem of reward generation to the task of generating 042 rough measures of *task progress*. Our framework uses LLMs to generate code for *progress functions*: 043 task-specific functions that map environment states to scalar measures of progress. For a given genre 044 of tasks, we follow the example from prior work on reward code generation (Yu et al., 2023) by asking practitioners to provide a small helper function library (ex. dist(x,y)) for the particular genre's observation spaces. Then, given both the helper function library and a single-sentence description of 046 a task within the domain (ex. "This environment require a closed door to be opened and the door 047 can only be pushed outward or initially open inward."), we leverage LLMs to generate the progress 048 functions. 049

We find that progress functions are most empirically effective when used within a count-based intrinsic reward framework: we treat the outputs of the progress functions as a simplified state space that groups together similar states from the environment, we discretize the states, and we compute state visitation counts across the discretized state space. Then, we treat the inverse square root of the visitation count as a count-based intrinsic reward, as in prior work (Kolter & Ng, 2009;

054 Tang et al., 2017), and learn policies using these count-based intrinsic rewards as task rewards. It 055 may seem straightforward to instead directly sum the outputs of the progress function and use the 056 sum as a reward – providing larger rewards for reaching states corresponding with greater progress 057 through the task. However, this approach neglects the typical reward scaling and weighting issues 058 common to dense reward shaping approaches (Sutton & Barto, 2018; Booth et al., 2023). Our progress-and-counts algorithm, ProgressCounts, achieves SOTA performance on the challenging Bi-DexHands benchmark, outperforming Eureka (Ma et al., 2023) by 4%. By limiting the role of 060 LLMs to generating progress functions and applying count-based intrinsic rewards to simplified 061 progress-based states, ProgressCounts achieves significantly greater sample efficiency compared to 062 approaches that use LLMs directly for estimating reward weighting and scaling (Yu et al., 2023; Ma 063 et al., 2023). Specifically, ProgressCounts requires 20 times fewer training runs than Eureka. 064

065 Specifically, we make the following contributions:

- 1. We re-frame the problem of reward generation in terms of generating coarse measures of *task progress*. Given a single-sentence task description and a small domain-specific library of helper functions, we leverage LLMs to generate *progress functions*.
- 2. We generate rewards from progress functions by treating the output of the progress function as a reduced state representation, discretizing progress in order to measure state visitation counts, and generating count-based intrinsic rewards.
- 3. We demonstrate that our algorithm ProgressCounts is effective by obtaining state-of-the-art performance on the Bi-DexHands benchmark. We outperform the prior state-of-the-art Eureka (Ma et al., 2023) by 4% while requiring $20 \times$ fewer samples.

2 RELATED WORK

079 080

066

067

068

069 070

071

073

075

076 077 078

Automated reward engineering. Learning directly from sparse rewards can be challenging (Ng et al., 1999; Hare, 2019; Vecerik et al., 2017). It is common for practitioners to carefully engineer dense reward functions to shape the learning process (Ng et al., 1999)—a labor-intensive (Sutton & Barto, 2018) and brittle (Booth et al., 2023) process. Advancements in automation such as Population Based Training (Jaderberg et al., 2017) have shown promise in refining this process by automating searches over fixed design spaces.

Foundation models offer the opportunity to automate reward engineering with powerful priors. The 087 output of foundation models can be used to propose tasks for open-ended learning curricula (Du 088 et al., 2023; Zhang et al., 2023) and add high-level structure to learning (Mirchandani et al., 2021). 089 Foundation models can also be used directly as rewards (Fan et al., 2022; Kwon et al., 2023; Sontakke 090 et al., 2024), and have the potential to generate reward function code, either scaffolded by reward 091 function templates (Yu et al., 2023) or in more freeform fashion (Ma et al., 2023; Venuto et al., 2024; 092 Li et al., 2024). We also leverage LLMs to inject task knowledge via code, but rather than attempting to generate complex reward functions, we simply ask them to identify a few key features associated 094 with task progress.

095 096

Count-based intrinsic rewards. Our algorithm leverages the idea of count-based intrinsic rewards in order to convert coarse progress-based state representations into rewards for policy learning. 098 Count-based intrinsic rewards are one of the main approaches to intrinsic motivation: estimating the "novelty" of a given state via state visitation counts (Tang et al., 2017). Approaches that 100 hash continuous spaces to discrete representations have shown considerable promise when the 101 discretization function is domain-specific (Tang et al., 2017; Ecoffet et al., 2021). In particular, the 102 Go-Explore algorithm (Ecoffet et al., 2021) pairs count-based intrinsic rewards along with simulator 103 state resets in order to achieve state-of-the-art results on several challenging Atari (Bellemare et al., 104 2013) tasks. The main downside to these approaches is that domain-aware discretization typically 105 requires significant human engineering (Tang et al., 2017; Ecoffet et al., 2021). In our algorithm, by discretizing task progress, we already have access to automated domain-specific state discretizations. 106 Absent the need for extensive human-engineered discretization functions, count-based intrinsic 107 rewards are both practical and effective for learning policies from progress functions.



Figure 1: ProgressCounts: an algorithm for reward generation via LLM-generated task 136 progress functions and count-based rewards. (A) We leverage a LLM to generate code for a 137 progress function, which distills task-specific features from a high-dimensional state space into a 138 low-dimensional notion of task progress. The LLM takes as input a high-level task description, a small library of feature engineering functions, and a description of the environment state space. On a 139 per-task basis, the user only needs to provide the task description as input. (B) We use heuristics to 140 discretize the output of the LLM-generated progress function, compute state visitation counts across 141 the discretized bins, and leverage standard count-based rewards to learn RL policies. 142

3 PRELIMINARIES

143 144

145 146

147

151

153

154

156

157

159

We consider the problem of automated reward generation for a sparse-reward task. The task is defined as a Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the 148 action space, \mathcal{P} is the transition probability function, and \mathcal{R} is a sparse reward function that provides 149 little guidance to the agent. The goal is to learn a policy $\pi : S \to \Delta(A)$ that maximizes the expected 150 cumulative reward $J(\pi) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t) \middle| s_0, \pi\right]$, where s_t, a_t are the state and action at time t.152

We assume the availability of three potential inputs for reward engineering:

- 1. A description of the features available in the environment (Figure 1-A, grey). We provide a description in the form of code, similar to Ma et al. (2023); Yu et al. (2023); Singh et al. (2023)
- 2. A short task description (Figure 1-A, green), similar to Ma et al. (2023); Yu et al. (2023).
- 3. An environment feature engineering library, offering a palette of additional, higher-level features that are not task-specific but may be generally useful for solving tasks in a given 161 domain (Figure 1-A, green). This is identical to the type of feature library in Yu et al. (2023).

162 Note that for (1), many learning scenarios with real-world deployment goals involve training in 163 simulators with access to environment code (Lin et al., 2024). For (2) and (3), an experienced 164 practitioner can quickly create a small feature engineering library, and the cost of making this library 165 is amortized across many tasks in the same domain.

166 167 168

4 **METHODS**

170 In this section, we first introduce our algorithm for leveraging LLM domain knowledge to generate progress functions, which distill key features from a high-dimensional environment state space to a 171 coarse low-dimensional notion of progress in the task. Then, we outline how we use the generated 172 progress functions: we view progress as measure of state, and leverage count-based intrinsic rewards 173 for learning.

174 175 176

4.1 PROGRESS FUNCTIONS

177 178

4.1.1 PROGRESS FUNCTION DEFINITION

179 Given a new task description, the first step of our process is to generate a progress function $P: S \to S$ R^k , which takes environment features $s \in S$ as input, and outputs information about the current 181 progress of the agent on the task. Especially for more complex tasks, it may be difficult to distill a 182 task to a single feature that tracks overall progress. Therefore, a progress function, given a state, is 183 asked to emit a positive scalar measure of progress for one or more subtasks. For instance, for the 184 SwingCup task (Figure 1-A), which involves 1) gripping the handles of the cup, 2) rotating the cup to 185 the correct orientation, a good progress function would break the task into two sub-tasks and return scalars measuring progress for both sub-tasks.

187 Specifically, the progress function outputs $[x_1, x_2, ..., x_k]$ where $x_i \in \mathcal{R}$ tracks task progress for 188 sub-task i. It also outputs additional variables $[y_1, y_2, ..., y_k]$ that inform our framework whether the 189 progress variables x_i are increasing or decreasing.

190 191

192

199

200

201

202

203

4.1.2 PROGRESS FUNCTION GENERATION

For any given task, domain knowledge is required in order to determine what features from the 193 environment are useful for assessing progress, and how to compute progress from those features. We 194 derive that domain knowledge from an LLM, which is used to generate code for the progress function 195 P. In order to translate domain knowledge into an effective progress function, we provide the LLM 196 with three inputs: 197

- 1. Function inputs: we specify the features available as inputs to the progress function via a description of the features in the environment state (Figure 1-A, grey). This information is available via the simulator.
- 2. Function outputs: we specify the desired output of the progress function via a short task description (Figure 1-A, green). Humans specify this information on a per-task basis.
- 204 3. Function logic: we structure the process of translating feature inputs to progress output by 205 providing the LLM with access to a environment feature engineering library (Figure 1-A, 206 green). This library offers a palette of additional, higher-level features to optionally use to 207 compute progress, and also indirectly suggests that certain types of feature transformations are beneficial to compute progress (ex. $l_2 dist(x, tgt)$). Humans create this library once per 208 genre or benchmark of tasks. 209
- 210 Note that on a per-task basis, ProgressCounts only requires a user to provide the short task description. 211 Please see Appendix A.3.2 for the libraries used for our Section 5 benchmarks. 212
- 213 Given the high-level task description, the small feature engineering library, and code describing the
- environment state space, we follow standard LLM prompting strategies to generate the progress 214 function code (Figure 1-A), blue). Please see Appendix A.4 for several examples of generated 215 progress functions.

216 4.2 FROM PROGRESS TO REWARD217

Given a perfect estimate of task progress, it may seem natural to directly use progress as a dense reward: for a given state s, compute the sum of the progress outputs $\mathcal{R}_{sum} = \sum_i x_i$, and use progress sum \mathcal{R}_{sum} as the reward for reaching s. However, learning from dense rewards can be brittle-small mistakes in reward design can often lead to a failure to learn effective policies (Booth et al., 2023). Progress functions offer highly simplified state representations—and given the coarse nature of these representations, we look for a more forgiving mechanism to generate rewards from these simplified representations. Therefore, we use a count-based intrinsic reward approach inspired by prior work that achieves state-of-the-art performance with domain-specific discretizations (Ecoffet et al., 2021).

233

234

235 236

237

243

244 245

246

247

248 249

250

253

4.2.1 PRELIMINARIES: COUNT-BASED REWARDS

To facilitate exploration, we leverage count-based rewards proportional to state novelty (Kolter & Ng, 2009): $n(s) \propto \frac{1}{\sqrt{c(s)}}$, where c(s) is the state visitation count. High-dimensional state spaces require a binning function $B : S \to S'$ that maps state space S to a (small) discrete space S' where we can tractably compute state visitation frequencies and novelty $n(s) \propto \frac{1}{\sqrt{c(B(s))}}$.

When learning a policy, sparse extrinsic rewards are augmented with a standard intrinsic reward proportional to n(s) (Tang et al., 2017):

$$\mathcal{R}_{total}(s_t, a_t) = \mathcal{R}(s_t, a_t) + \lambda_c n(B(s_{t+1})) \tag{1}$$

Where s_t is the state at time t, a_t is the action, s_{t+1} is the next state, and λ_c is a hyperparameter weighting the intrinsic reward relative to extrinsic reward $\mathcal{R}(s_t, a_t)$.

Effective state binning should encode information if and only if it is relevant to solving the particular task (Tang et al., 2017; Ecoffet et al., 2021).

4.2.2 COUNT-BASED REWARDS FROM PROGRESS

We automatically generate a task-specific binning function B by converting the continuous progress values P(s) emitted by the progress function to discrete states using a mapping $D: R^k \to S'$ that:

- 1. Estimates relevant value ranges (min_i, max_i) for each x_i from progress data
- 2. Discretizes within (min_i, max_i) to produce discrete progress features x'_i . We discretize later subtasks with finer granularity in order to encourage more exploration closer to the goal.
- 3. Defines $B(s) = D(P(s)) = \sum_{i} x'_{i}$.

Heuristic discretization avoids the need to learn scalar progress ranges from environment interaction, an approach used in prior work (Shinn et al., 2024; Ma et al., 2023). While we could have leveraged the LLM directly to emit logic for discretizing progress features, we chose to use these heuristics instead since LLMs are known to struggle with numerical reasoning (Shen et al., 2023). Details and example discretization code are included in Appendix A.6.

Having defined mapping *B* to discretize progress features, we measure state novelty from bin visitation counts via $n(s) \propto \frac{1}{\sqrt{c(B(s))}}$, augmenting existing sparse extrinsic rewards. As shown in Figure 1-B, we learn policies using Proximal Policy Optimization (PPO) (Schulman et al., 2017), augmenting the sparse extrinsic task rewards with intrinsic rewards via count-based intrinsic motivation (See Eq. 1).

264 265

5 EVALUATION

266 267

We evaluate ProgressCounts by using it to train policies on Bi-DexHands: a challenging sparse-reward
 benchmark consisting of 20 bimanual manipulation tasks. We also include additional results from the MiniGrid benchmark in the Appendix.



Figure 2: On the Bi-DexHands benchmark, ProgressCounts produces policies that perform comparably to those of Eureka in terms of average task success rate, at a much smaller sample 285 **budget.** Eureka's evolutionary algorithm requires 48 policy samples (training runs with different 286 generated reward functions) to find a policy whose performance matches that of human-designed dense reward functions. ProgressCounts requires only four policy samples (different progress functions), generating a policy that outperforms the human-designed baseline and exceeding the peak performance achieved by Eureka after 80 policy samples ($20 \times$ the cost of ProgressCounts). 289

290 291

292

270

271

272

273 274 275

276 277 278

279

281

282 283

284

287

288

5.1 EXPERIMENTAL SETUP

293 Bi-DexHands. The Bi-Dexterous Manipulation benchmark (Chen et al., 2022) (Bi-DexHands) consists of 20 bimanual manipulation tasks with continuous state and action spaces, such as using two 295 robotic hands to lift a pot or simultaneously pass objects between the hands. These tasks have sparse 296 rewards and require complex coordinated motion, making them a challenging test for leveraging 297 language models to guide policy learning. Following conventions from prior work (Ma et al., 2023), 298 progress functions have acess to the environment state space code, and we evaluate performance on 299 Bi-DexHands in terms of the policy's success rate at completing each task, averaged over five trials (policy training runs with different seeds). 300

301 We use Bi-DexHands to evaluate the policy performance and sample efficiency of ProgressCounts 302 against three baselines. 1) **Sparse** extrinsic rewards upon task success, 2) **Dense**: expert-written dense 303 extrinsic rewards from the original benchmark, 3) rewards generated using the Eureka LLM-based 304 reward generation algorithm (Ma et al., 2023), the current state-of-the-art reward generation method on Bi-DexHands. 305

306

307 **Training configuration.** In all experiments we train policies using PPO (Schulman et al., 2017). We train policies using the PPO hyperparameters and sample budgets (100M environment samples) 308 established by Bi-DexHands Chen et al. (2022), also used in prior work Eureka Ma et al. (2023). We 309 set the intrinsic reward coefficient $\lambda_c = 1e - 3$, and discretize progress into 1000 bins. We leverage 310 GPT-4-Turbo ('gpt-4-turbo-2024-04-09') as the LLM (Achiam et al., 2023) used to generate progress 311 functions. Following the experimental procedure from prior work (Ma et al., 2023), we use the 312 LLM to generate multiple options for the progress function, and select the resulting policy that 313 achieves the highest success from a single training run-we refer to the different trained policies 314 as *policy samples*. Unless otherwise specified, ProgressCounts uses four policy samples per task, 315 and all policies are trained using 100M environment samples (number of environments \times number of 316 simulation steps). All LLM prompts, including task descriptions and environment feature engineering 317 primitives, are included in Appendix A.3.

318

319 5.2 COMPARISON TO EXTRINSIC REWARD BASELINES 320

321 ProgressCounts trains policies that (on average) outperform those from Eureka on Bi-DexHands, using only 5% of Eureka's training budget. Averaged over all Bi-DexHands tasks, 322 ProgressCounts achieves a success rate of 0.59, 13% higher than human-written dense rewards, and 323 4% higher than Eureka, the state-of-the-art method on this benchmark (Figure 2). Most importantly,



Bi-DexHands tasks. See Appendix A.7 for results in tabular form.

TwoCatchUnderarm Success Rate (LLMCount: Single Policy Sample)



Figure 4: By allocating many environment samples to a single training run, ProgressCounts trains a policy that achieved high success on the challenging TwoCatchUnderarm task. All baselines achieved zero success on this task given a two billion environment sample budget.

Eureka's evolutionary algorithm requires 80 policy samples (generated reward functions) to find good reward functions for the task. In contrast, ProgressCounts only requires four policy samples (generated progress functions). By structuring reward engineering around a constrained progress function and heuristic discretization, we reduce the unreliability associated with using LLMs for unconstrained code generation Yu et al. (2023). This approach also allows for more robust state discretization for count-based intrinsic rewards, using only a limited number of progress function generation attempts, *without requiring costly feedback-driven evolution*.

Figure 3 presents policy performance for all 20 tasks in the Bi-DexHands benchmark. Across the
 benchmark, ProgressCounts matches or exceeds Eureka in performance on 13 of the 20 tasks, and
 ProgressCounts matches or exceeds the performance of the expert-written dense reward (Human) on
 17 of the 20 tasks.

369

339

340 341

342

343

349 350

351

352 353 354

355

356 357

370 Given the same environment sample budget as Eureka, ProgressCounts can produce higher-371 **performance policies.** Since ProgressCounts requires fewer policy samples to find good policies, 372 users can more confidently allocate significant fractions of a training budget to a small number of 373 policy samples. We use ProgressCounts to train four policies on the TwoCatchUnderarm task, each 374 for a total of two billion environment samples (the same number of samples used in aggregate, across 375 all policy samples, for Eureka training). The best resulting policy achieves a task success rate of 0.55, and continues to improve with further samples (Figure 4). On the other hand, all extrinsic 376 reward baselines, as well as ProgressCounts trained on 400 million environment samples (four policy 377 samples, with 100 million environment samples each), achieve a success rate of nearly zero on this

task (Figure 3). To our knowledge, ProgressCounts is the first method to achieve reasonable success on this challenging task.

5.3 METHOD ABLATIONS

384	Task name	ProgressCounts	ProgressAsReward	SimHashCounts
385	Average	0.59	0.45	0.34
386	Over	0.93	0.90	0.91
387	DoorCloseInward	1.00	1.00	1.00
388	DoorCloseOutward	0.90	1.00	0.76
389	DoorOpenInward	0.07	0.00	0.00
390	DoorOpenOutward	0.99	0.31	0.99
391	Scissors	1.00	1.00	1.00
392	SwingCup	0.97	0.99	0.94
393	Switch	0.00	0.00	0.00
394	Kettle	0.83	0.00	0.00
395	LiftUnderarm	0.22	0.08	0.00
396	Pen	0.49	0.22	0.09
397	BottleCap	0.94	0.04	0.94
202	CatchAbreast	0.56	0.49	0.00
200	CatchOver2UnderArm	0.90	0.94	0.00
399	CatchUnderarm	0.76	0.88	0.00
400	ReOrientation	0.03	0.06	0.02
401	GraspAndPlace	0.99	0.98	0.08
402	BlockStack	0.05	0.00	0.06
403	PushBlock	0.03	0.02	0.01
404	TwoCatchUnderarm	0.03	0.01	0.00

405

381

382

406 Table 1: An ablation testing whether Progress Functions and Count-Based Rewards are both 407 necessary for ProgressCounts across the 20 tasks in Bi-DexHands. ProgressCounts is our 408 algorithm. ProgressAsReward takes the best generated progress functions, and directly uses the 409 summed progress variables as a dense reward function. SimHashCounts applies SimHash to the observation space as the binning function instead of progress-based bins (the method from [4]). 410 Results are averaged across 5 trials for **ProgressCounts**, and are single-trial numbers for the ablated 411 methods. ProgressCounts requires both key components of the algorithm for an 0.59 average 412 task success rate. 413

414

The success of ProgressCounts is due to both the use of progress functions to collapse simulator states into bins and due to the effectiveness of count-based intrinsic exploration applied to these bins. While progress function might seem a suitable dense reward, progress-based rewards only achieve a success rate of 0.45 (Table 1), so best performance is achieved when using count-based exploration across discretized progress bins. Both LLM-generated progress functions and count-based intrinsic exploration are necessary to achieve our SOTA performance.

420 421

Progress functions generate more effective bins than SimHash: On Bi-DexHands, Table 1 high-lights that ProgressCounts achieves a success rate of 0.59 with progress-based bins, and only achieves a success rate of 0.34 with SimHash-based (Sadowski & Levin, 2007) bins across the observation space (Tang et al., 2017). Across the benchmark, progress-based bins achieve performance equal to or better than SimHash-based bins across 19 of 20 tasks, with the remaining task (BlockStack) within the margin of error (see Table 8 in the Appendix for standard deviations). This result aligns with prior work on human-written hash functions for count-based rewards, where the integration of domain knowledge consistently improves performance (Tang et al., 2017; Ecoffet et al., 2021).

429

430 **Count-based rewards are more effective than directly using progress as reward** While the 431 sum of the outputs of a progress function $\mathcal{R}_{sum} = \sum_i x_i$ might seem a viable reward signal for learning, Table 1 illustrates that progress-based dense rewards only achieve a success rate of 0.45,

2	Task	Default	No feature library	No heuristic discretization
4	SwingCup	0.97	0.90	0.00
5	CatchUnderarm	0.76	0.00	0.76
36	DoorCloseOutward	0.90	0.86	0.92

Table 2: Both the environment feature library and heuristic progress discretization help task success rate. When the feature engineering library is removed from the LLM prompt, we obtain comparable performance on SwingCup and DoorCloseOutward, but the CatchUnderarm policy completely fails to learn. When we ask the LLM to directly generate code for discrete bins (removing heuristic discretization), SwingCup fails to learn an effective policy.

while progress functions paired with count-based intrinsic rewards achieve a success rate of 0.59, matching or outperforming the dense-reward alternative on 15 of 20 tasks. This result highlights that, given coarse state representations from progress functions, count-based intrinsic rewards are more effective than standard dense rewards.

449 The environment feature engineering library helps generate effective progress features In Table 2, we ablate the impact of providing a feature engineering library to help ProgressCounts with 450 generating code to compute progress features. When we remove the feature library, we still obtain 451 comparable performance on both SwingCup (task success rate of 0.97 vs. 0.90) and DoorCloseOut-452 ward (0.90 vs. 0.86), but the CatchUnderarm policy fails to learn completely. Upon inspection of 453 the generated code in Appendix A.5, the LLM chooses to incorporate object linear velocity into the 454 progress function, a variable that is not directly relevant to task success. This variable is averaged 455 with more relevant variables to derive an overall progress metric. Having incorrectly modeled task 456 progress, the policy's unnecessary exploration is likely responsible for task failure. The library for 457 Bi-DexHands (included in Appendix A.5) only contains functions measuring Euclidean and rotational 458 distance; we hypothesize that knowledge of transformations available in the feature library helps the 459 LLM ignore features for which the library is not applicable (ex. the irrelevant velocity features).

- 461 **ProgressCounts benefits from using heuristics to discretize progress features** In Table 2, we 462 also ablate the impact of using heuristics to discretize and combine progress features-instead, 463 we ask the LLM to directly generate code to output discrete bins corresponding to task progress. We obtain comparable task success rate on DoorCloseOutward (0.90 vs. 0.92) and CatchUnderarm 464 (0.76 vs. 0.76), but the trained policies completely fail without heuristic discretization for SwingCup-465 as seen in Appendix A.5, the LLM incorrectly guesses the relevant range of values for multiple 466 features, and as a result the binning function is ineffective for facilitating task-relevant exploration. 467 Heuristic discretization helps avert this failure mode when constructing task-specific state binning 468 functions. 469
- 470

460

437 438

439 440

441

442 443 444

445

446

447

448

6 DISCUSSION

471 472 473

The state-of-the-art results achieved by ProgressCounts demonstrate two key takeaways:

474 First, LLM-generated progress functions offer a compelling mechanism to generate coarse task-475 specific state representations, and alongside count-based intrinsic rewards offer an empirically-476 superior alternative to using LLMs to engineer reward functions. ProgressCounts outperforms Eureka, 477 which uses LLMs to engineer reward functions, both in terms of performance and sample efficiency. One reason for this success is the structure (task progress, count-based rewards, etc.) we build into 478 the ProgressCounts framework, which increases the quality and reliability of LLM responses, and 479 reduces the need for trial and error across reward weights and scaling. Perhaps more interestingly, 480 we hypothesize that ProgressCounts also benefits from count-based intrinsic rewards being robust to 481 non-optimal binning functions, unlike reward functions where even minor errors can easily lead to a 482 failure to solve tasks successfully. 483

484 Second, despite being relatively under-utilized in recent research, count-based intrinsic rewards can
 485 be surprisingly effective at training policies that operate in complex high-dimensional state spaces
 when given an adequate binning function. Interestingly, the results achieved by ProgressCounts

486 suggest that these binning functions do not need to be hugely complex; ProgressCounts outperforms 487 state-of-the-art, human-engineered dense reward functions using count-based exploration driven by 488 binning functions that contain less than 20 lines of code.

489 Overall, we believe ProgressCounts represents a novel and promising strategy for injecting domain 490 knowledge from large language models into an RL training loop. We hope that these results will 491 encourage further research into (and more general usage of) count-based intrinsic methods, as well as 492 exploration of other novel methods for leveraging LLMs to assist in solving reinforcement learning 493 tasks.

494 495 496

497

498

499

500

521

522

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environ-501 ment: An evaluation platform for general agents. Journal of Artificial Intelligence Research, 47: 502 253-279, 2013.
- 504 Serena Booth, W Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi. The 505 perils of trial-and-error reward design: misdesign through overfitting and invalid task specifications. 506 In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 5920–5929, 2023. 507
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network 508 distillation. arXiv preprint arXiv:1810.12894, 2018. 509
- 510 Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong Feng, Jiechuan Jiang, Zongqing Lu, Stephen 511 McAleer, Hao Dong, Song-Chun Zhu, and Yaodong Yang. Towards human-level bimanual 512 dexterous manipulation with reinforcement learning. Advances in Neural Information Processing 513 Systems, 35:5150-5163, 2022. 514
- 515 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo Perez-Vicente, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & 516 customizable reinforcement learning environments for goal-oriented tasks. Advances in Neural 517 Information Processing Systems, 36, 2024. 518
- 519 Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek 520 Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. In International Conference on Machine Learning, pp. 8657-8677. PMLR, 2023.
- 523 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then explore. Nature, 590(7847):580-586, 2021. 524
- 525 Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, 526 De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied 527 agents with internet-scale knowledge. Advances in Neural Information Processing Systems, 35: 528 18343-18362, 2022. 529
- 530 Joshua Hare. Dealing with sparse rewards in reinforcement learning. arXiv preprint arXiv:1910.09281, 2019. 531
- 532 Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali 533 Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. Population based training 534 of neural networks. arXiv preprint arXiv:1711.09846, 2017. 535
- J Zico Kolter and Andrew Y Ng. Near-bayesian exploration in polynomial time. In Proceedings of the 26th annual international conference on machine learning, pp. 513–520, 2009. 538
- Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language 539 models. arXiv preprint arXiv:2303.00001, 2023.

540 541 542 543	Hao Li, Xue Yang, Zhaokai Wang, Xizhou Zhu, Jie Zhou, Yu Qiao, Xiaogang Wang, Hongsheng Li, Lewei Lu, and Jifeng Dai. Auto mc-reward: Automated dense reward design with large language models for minecraft. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and</i> <i>Pattern Recognition</i> , pp. 16426–16435, 2024.
544 545 546	Toru Lin, Zhao-Heng Yin, Haozhi Qi, Pieter Abbeel, and Jitendra Malik. Twisting lids off with two hands. <i>arXiv preprint arXiv:2403.02338</i> , 2024.
547 548 549	Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. <i>arXiv preprint arXiv:2310.12931</i> , 2023.
550 551	Denys Makoviichuk and Viktor Makoviychuk. rl-games: A high-performance framework for rein- forcement learning. https://github.com/Denys88/rl_games, May 2021.
552 553 554	Suvir Mirchandani, Siddharth Karamcheti, and Dorsa Sadigh. Ella: Exploration through learned language abstraction. <i>Advances in neural information processing systems</i> , 34:29529–29540, 2021.
555 556	Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In <i>Icml</i> , volume 99, pp. 278–287, 1999.
557	Caitlin Sadowski and Greg Levin. Simhash: Hash-based similarity detection, 2007.
558 559 560	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
561 562	Ruoqi Shen, Sébastien Bubeck, Ronen Eldan, Yin Tat Lee, Yuanzhi Li, and Yi Zhang. Positional description matters for transformers arithmetic. <i>arXiv preprint arXiv:2311.14737</i> , 2023.
563 564 565	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 36, 2024.
566 567 568 569	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 11523–11530. IEEE, 2023.
570 571 572 573	Sumedh Sontakke, Jesse Zhang, Séb Arnold, Karl Pertsch, Erdem Bıyık, Dorsa Sadigh, Chelsea Finn, and Laurent Itti. Roboclip: One demonstration is enough to learn robot policies. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
574	Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
575 576 577 578	Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, OpenAI Xi Chen, Yan Duan, John Schulman, Filip DeTurck, and Pieter Abbeel. # exploration: A study of count-based exploration for deep reinforcement learning. <i>Advances in neural information processing systems</i> , 30, 2017.
579 580 581 582	Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. <i>arXiv preprint arXiv:1707.08817</i> , 2017.
583 584 585	David Venuto, Sami Nur Islam, Martin Klissarov, Doina Precup, Sherry Yang, and Ankit Anand. Code as reward: Empowering reinforcement learning with vlms. <i>arXiv preprint arXiv:2402.04764</i> , 2024.
586 587 588	Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language to rewards for robotic skill synthesis. <i>arXiv preprint arXiv:2306.08647</i> , 2023.
589 590 591	Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. Omni: Open-endedness via models of human notions of interestingness. <i>arXiv preprint arXiv:2306.01711</i> , 2023.
592 593	Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E Gonzalez, and Yuandong Tian. Noveld: A simple yet effective exploration criterion. <i>Advances in Neural Information Processing Systems</i> , 34:25217–25230, 2021.

594 A APPENDIX

596

597

624

625

626 627

628

A.1 ABLATING THE CHOICE OF CONSTANTS FOR PROGRESSCOUNTS

598 ProgressCounts does not require search over hyperparameters on a per-task basis-we set the hyperparameters once per benchmark and did not tune them (following the same experimental protocol as Eureka (Ma et al., 2023)). We evaluate ProgressCounts's robustness to different hyperparameters 600 using two ablations: 1) the intrinsic reward weight λ_c , and 2) the number of discrete bins used, each 601 tested across three tasks from Bi-DexHands. Table 3 shows similar task performance for λ_c of 1e-2, 602 1e-3, and 1e-4, with the exception of slightly lower task performance on CatchUnderarm for 603 1e-4. Table 4 shows similar task performance with 500, 1000, and 2000 bins across the three tasks. 604 There are no clear trends in performance across parameters in either table, and our chosen parameters 605 $(\lambda_c = 1e - 3, 1000 \text{ bins})$ are actually sub-optimal for most tasks. 606

Task	$\lambda_c = 0.01$	$\lambda_c = 0.001$	$\lambda_c = 0.0001$
SwingCup	0.95	0.97	0.98
CatchUnderarm	0.8	0.76	0.4
DoorCloseOutward	0.99	0.9	0.85

Table 3: Ablating the impact of intrinsic reward coefficient λ_c on ProgressCounts performance. We report success rates for 3 values across 3 Bi-DexHands tasks. Success rates are averaged across 5 trials. Default is $\lambda_c = 0.001$

Task	500 bins	1000 bins	2000 bins
SwingCup	0.97	0.97	0.97
CatchUnderarm	0.78	0.76	0.67
DoorCloseOutward	0.84	0.9	0.99

Table 4: Ablating the impact of the number of bins on ProgressCounts performance. We report success rates for 3 values across 3 Bi-DexHands tasks. Success rates are averaged across 5 trials. Default is 1000 bins.

A.2 MINIGRID EXPERIMENTS

Having demonstrated that ProgressCounts yields state-of-the-art performance on the Bi-DexHands
benchmark, we further evaluate how well the progress-based counts from ProgressCounts can serve
as a novelty metric within more sophisticated intrinsic motivation algorithms. Specifically, we test
how well progress-based novelty from Section 4.2 performs as a novelty measure within the NovelD
meta-criterion (Zhang et al., 2021), which provides reward proportional to the difference in novelty
between consecutive states:

635 636

637

$$\mathcal{R}_{total}(s_t, a_t) = \mathcal{R}(s_t, a_t) + \lambda_c \max(n(B(s_{t+1})) - \alpha n(B(s_t)), 0) \mathbb{1}[n_e(s_{t+1}) = 1]$$
(2)

638 Where $\alpha \in [0, 1]$ discounts previous novelty, and $n_e(s_{t+1})$ measures episodic novelty, measuring the 639 number of visits to a state within the current episode.

640 Using the MiniGrid benchmark (Chevalier-Boisvert et al., 2024), we evaluate on a subset of eight 641 difficult exploration tasks across two task distributions: four KeyCorridor variants and four Ob-642 structedMaze variants. These tasks provide sparse rewards upon task success, and these rewards 643 are proportional to the efficiency of completing the goal. We compare the efficacy of the NovelD 644 exploration meta-criterion (Zhang et al., 2021) when measuring novelty with **ProgressCounts** as 645 well as **RND** Burda et al. (2018), the method used in the original NovelD paper. Following Zhang et al. (2021), we measure episode rewards averaged over four trials when combining the novelty 646 metrics with the NovelD algorithm-we measure performance in terms of the samples required to 647 reach a threshold task reward. We set $\lambda_c = 0.1$, and we do not require progress discretization on this

Env Type	Layout	ProgressCounts	RND
KeyCorridor (Medium)	S3R3	0.2	0.5
	S4R3	0.5	0.9
	S5R3	0.9	1.3
	S6R3	1.2	1.7
ObstructedMaze (Hard)	2Dlhb	1.6	4.0
	1Q	1.0	2.8
	2Q	2.5	4.7
	Full	2.9	7.2

Table 5: ProgressCounts is more sample-efficient than RND when used as a novelty metric within NovelD. We measure the number of environment samples ($\times 10^7$) required for NovelD to pass a threshold reward of 0.75 (except for ObstructedMaze-Full, where computational constraints limit us to a threshold of 0.5). Across four variants of the KeyCorridor task and four variants of the ObstructedMaze task, ProgressCounts is up to 64% more sample-efficient than RND.

environment since the progress functions are already discrete. As with the Bi-DexHands experiments, we run 4 trials of progress function generation.

Table 5 compares the performance of ProgressCounts and RND as a novelty metric within NovelD on MiniGrid. Across both KeyCorridor and ObstructedMaze families of tasks, ProgressCounts improves the sample efficiency of NovelD at reaching a threshold reward compared to RND. This trend holds across progressively more complicated tasks: ProgressCounts improves sample efficiency by 60% for KeyCorridorS3R3, the simplest task, and also improves sample efficiency by 60% on ObstructedMaze-Full, the hardest task. Note that ProgressCounts is also more computationally efficient than RND, which learns an additional network to output intrinsic rewards (Burda et al., 2018). Full training curves are in Appendix A.9.

A.3 DETAILS ON LLM INPUTS

A.3.1 SYSTEM PROMPT

681	1	You are a reinforcement learning engineer trying to write progress
682		
683		→ as possible.
684	2	Your goal is to identify the variables for the environment that are
685		→ maximally relevant for measuring progress in the task described
686	3	Some tasks may have only a single stage, and some tasks may have two
687	-	→ separate stages.
688	4	You will be provided with a definition of the observation space for
689		\hookrightarrow a reinforcement learning environment, and also provided with a
690		\rightarrow small set of helper functions that can be used to transform the
691	5	\rightarrow variables in the observation space.
692	5	while a function that feturis the variable most associated with task \rightarrow progress for each stage of the task.
693	6	This function can take as input any member of self defined in
69/		↔ compute_observations, and can apply any of the helper functions
605		\hookrightarrow to any variables from self.obs_buf to generate new derived
606		\rightarrow features (ex. computing the distance between object and goal).
090		→ If a single stage requires multiple progress variables, average
697	7	\rightarrow the valuables. Also return a bool for each variable that is True if progress
698	,	\rightarrow requires the variable to increase, and False if it requires the
699		\rightarrow variable to decrease.
700	8	
701	9	Function signature:

A.3.2 ENVIRONMENT FEATURE ENGINEERING LIBRARY

Bi-DexHands For the Bi-DexHands benchmark, our library is composed of three simple functions:
1) Euclidean distance, 2) rotational distance between quaternions, 3) Euclidean distance to a 'goal' state if it exists.

```
1 # Determine distance of an object from a "goal", if it exists
2 def goal_dist(self, x):
       return torch.norm(self.goal_pos - x, p=2, dim=-1)
3
5 # Determine distance between two objects
6 def dist(self, x, y):
      return torch.norm(x - y, p=2, dim=-1)
7
8
9 # Rotational distance
10 def rot_dist(self, object_rot, target_rot):
      quat_diff = quat_mul(object_rot, quat_conjugate(target_rot))
11
      rot_dist = 2.0 * torch.asin(torch.clamp(torch.norm(quat_diff[:,
12
       \leftrightarrow 0:3], p=2, dim=-1), max=1.0))
      return rot_dist
13
```

MiniGrid For the MiniGrid benchmark, our library is composed of three simple functions: 1) breadth-first search to find the shortest path between two grid cells, accounting for walls, 2) a function finding the grid position of a given type of object, 3) a function that finds the grid position of an object, given that the object is on the path between two given locations.

```
731
          1 def bfs(grid, start, end):
         2
732
                Perform BFS to find the shortest path from start to end in a
         3
733
                → minigrid environment.
734
         4
                Args:
735
                 - grid (np.array): The grid represented as a numpy array of
         5
736
                 \hookrightarrow shape (n, m, 3).
                - start (tuple): Starting position (x, y).
          6
737
                - end (tuple): Ending position (x, y).
         7
738
         8
739
                Returns:
         9
740
                - path (list): List of tuples as coordinates for the shortest
         10
                 \rightarrow path, including start and end.
741
                                  Returns an empty list if no path is found.
         11
742
                .....
         12
743
         13
                queue = deque([start])
744
                paths = {start: [start]}
         14
745
                directions = [(1, 0), (0, 1), (-1, 0), (0, -1)] # Down, right,
         15
746
                 \leftrightarrow up, left
                while queue:
         16
747
         17
                     current = queue.popleft()
748
                     #if grid[current[0], current[1], 0] != 1:
         18
749
                          print("Current", current)#, paths[current])
         19
                     #
750
                     #
                          print("Content", grid[current[0], current[1]])
         20
                     if current == end:
751
         21
         22
                         return paths[current]
752
                     for direction in directions:
         23
753
                         neighbor = (current[0] + direction[0], current[1] +
         24
754
                          \leftrightarrow direction[1])
755
                         if (0 <= neighbor[0] < grid.shape[0] and</pre>
         25
```

```
756
                              0 <= neighbor[1] < grid.shape[1] and</pre>
         26
757
                             neighbor not in paths and
         27
758
                             is_traversable(grid[neighbor])):
         28
759
         29
                             paths[neighbor] = paths[current] + [neighbor]
                             queue.append(neighbor)
         30
760
                return [] # Return an empty list if no path is found
         31
761
         32
762
         33 def get_position(grid, object_type, color=None):
763
         34
764
                 Get the position of the object of the specified type in the
         35
                 \hookrightarrow grid.
765
                Args:
         36
766
                 - grid (np.array): The grid represented as a numpy array of
         37
767
                 \rightarrow shape (n, m, 3).
768
                 - object_type (int): The type of the object to find.
         38
769
         39
                Returns:
770
         40
         41
                 - position (tuple): The position of the object in the grid.
771
                 42
772
                for i in range(grid.shape[0]):
         43
773
                     for j in range(grid.shape[1]):
         44
774
         45
                         if grid[i, j, 0] == object_type:
                              if color is not None and grid[i, j, 1] != color:
775
         46
                                  continue
         47
776
                              return (i, j)
         48
777
         49
                return None
778
         50
779
         51 def get_position_on_path(grid, agent_pos, final_pos, object_type,
            ↔ color=None, closed=None):
780
                path = bfs(grid, agent_pos, final_pos)
         52
781
         53
                for pos in path:
782
         54
                     if grid[pos[0], pos[1], 0] == object_type:
783
                         if color is not None and grid[pos[0], pos[1], 1] !=
         55
784
                             color:
                          \hookrightarrow
785
         56
                             continue
                         if closed is not None and grid[pos[0], pos[1], 2] !=
         57
786
                            closed:
787
         58
                              continue
788
                         return pos
         59
789
         60
                return None
790
791
       A.3.3 TASK DESCRIPTIONS
792
793
```

704	Bi-DexHands Environments
794	Task name
795	Task description
796	Task success condition
797	
798	Over
799	This environment requires an object in one hand to be thrown to the goal location on the other hand.
800	The task is a simple single-stage task.
801	1[dist < 0.03]
802	DoorCloseInward
803	This environment require a closed door to be opened and the door can only be pushed outward or
804	initially open inward.
805	$1[door_handle_dist < 0.5]$
806	DoorCloseOutward
807	This environment requires a closed door to be opened, but because they can't complete the task by
808	simply nushing we need to catch the handle by hand and then open it so it is relatively difficult
809	$1[door_handle_dist < 0.5]$

0	
1	DoorOpenInward
2	This environment requires the hands to grab the handles of the doors, then pull the two doors apart.
3	$1[door_handle_dist > 0.5]$
4	DoorOpenOutward
5	This environment requires the hands to grab the handles of the doors, then pull the two doors apart
6	$1[door_handle_dist < 0.5]$
7	Seiscone
8	Scissors This environment requires the hands to grab the handles of a pair of scissors, then open the scissors
9	1[dof nos > -0.3]
20	
21	SwingCup This anticomment involves two hands and a dual handle our, we need to use two hands to hold and
22	This environment involves two hands and a dual nancie cup, we need to use two hands to hold and swing the cup to a target orientation
23	swing the cup to a target orientation. 1[rot dist < 0.785]
24	
25	Switch
26	I his environment requires both hands to reach their respective switches, then lower the switch handle
27	positions by applying a strong downward force. 1[1.4] = (loft whiteh z + wight whiteh z) > 0.05]
.8	$1[1.4 - (left_Switch_2 + llght_Switch_2) > 0.05]$
29	Kettle
)	This environment requires the hands to grab the kettle, then move the kettle spout to the bucket.
	$1[\textit{pucket} - \textit{kettle_spout}] < 0.05]$
	LiftUnderarm
	This environment requires grasping the pot handle with two hands and lifting the pot to the designated
	position.
	1[dist < 0.05]
	Pen
	This environment requires the cap to be removed from the pen.
	$1[5 \times pen_cap - pen_body > 1.5]$
	BottleCap
	This environment involves two hands and a dual handle cup, we need to move the cap away from the
	bottle.
	1[dist > 0.03]
	CatchAbreast
	This environment consists of two shadow hands placed side by side in the same direction and an
	object that needs to be passed from one palm to a goal position on the other.
	1[dist < 0.03]
	CatchOver2Underarm
	This environment requires an object in one hand to be thrown to the goal location on the other hand
	1[dist < 0.03]
	CotabUndowarm
	Calculurianin This environment requires an object in one hand to be thrown to the goal location on the other hand
	1[dist < 0.03]
	ReOrientation
	This environment involves two hands and two objects. Each hand holds an object and we need to
	reorient the object to the target orientation.
	$1[rot_dist < 0.1]$
	GraspAndPlace
	This environment consists of dual-hands, an object and a bucket that requires us to pick up the object
	and put it into the bucket.
	1 block - bucket < 0.2]
	BlockStack

863 This environment involves dual hands and two blocks, and we need to stack the block as a tower.

PushBlock	
This enviror	ment involves dual hands and two blocks, and we need to push both blocks to goa
positions.	
$1 0.1 \le left_{-} \\ 1 0.5 \times left_{-} $	$dist \leq 0.1$ and $right_dist \leq 0.1$ $dist - 0.1$ and $right_dist \leq 0.1$]
TwoCatchU	nderarm
This environ $1[dist < 0.03]$	ment requires two objects to be thrown into the other hand at the same time. 3]
	MiniGrid Environments
	Task name
	Task description Sparse reward upon episode and
	Sparse reward upon episode end
KevCorrido	r
This environ	ment is a Gridworld that requires the agent to navigate to a key, then to a blue ball.
A reward of	$1 - 0.9 * (step_count/max_steps)$ is given for success, and 0 for failure.
Obstructed	Maze
This enviror	ment is a Gridworld that requires the agent to navigate to a blue ball, and there is
locked door	on the path, and multiple misleading locked doors in the environment. Do not worr
about the key	у.
A reward of	$1 - 0.9 * (step_count/max_steps)$ is given for success, and 0 for failure.

A.4 PROGRESS FUNCTION SAMPLES

SwingCup:

922	1	<pre>def progress_function(self) -> Tuple[List[torch.Tensor],</pre>
923		↔ List[bool]]:
924	2	<pre># Stage 1: Gripping Handles</pre>
925	3	# Calculate distances of hands to respective handles
926	4	<pre>dist_left_hand_to_left_handle = self.dist(self.left_hand_pos,</pre>
927		<pre> self.cup_left_handle_pos) </pre>
000	5	<pre>dist_right_hand_to_right_handle = self.dist(self.right_hand_pos,</pre>
920		Self.cup_right_handle_pos) # Average distance to bandles for gripping as one program.
929	0	# Average distance to nandles for gripping as one progress
930	7	\rightarrow metric
931		$avg_giip_uistance = (uist_ieit_inanu_co_ieit_inanuie + dist_right_hand to_right_handle) / 2$
932	8	/ dibe_right_hand_co_right_handic, / 2
933	9	# Stage 2: Orienting Cup
934	10	# Calculate rotational distance to the goal orientation
935	11	<pre>orientation_error = self.rot_dist(self.object_rot,</pre>
936		\leftrightarrow self.goal_rot)
937	12	
022	13	# collect progress variables and directions (irue for increasing,
930	14	\rightarrow <i>false for decreasing)</i>
939	14	progress_vals = [avg_gilp_ulstance, olientation_elion]
940	15	progress_arrections = [raise, raise] # Both Shourd decrease
941	17	return progress vars, progress directions
942		

CatchUnderarm

```
1 def progress_function(self):
      # Calculate distance between the object's current position and
2
      \rightarrow the goal.
      object_to_goal_distance = self.goal_dist(self.object_pos) #
3
       ↔ Calculated via the helper function.
4
      # Since the task progress is measured by how close the object
5
      \hookrightarrow comes to the goal,
      # we want this distance to decrease. Thus, progression in the
6
       → task is achieved
      # when `object_to_goal_distance` is minimized.
7
      progress_vars = [object_to_goal_distance]
8
      progress_directions = [False] # False indicates that we require
9
       \leftrightarrow the distance to decrease to indicate progress.
10
      return progress_vars, progress_directions
11
```

DoorCloseOutward

```
1 def progress_function(self):
      # Calculate distance between the object's current position and
2
       \leftrightarrow the goal.
      object_to_goal_distance = self.goal_dist(self.object_pos) #
3
       \leftrightarrow Calculated via the helper function.
4
      # Since the task progress is measured by how close the object
5
      \leftrightarrow comes to the goal,
      # we want this distance to decrease. Thus, progression in the
6
      → task is achieved
      # when `object_to_goal_distance` is minimized.
7
      progress_vars = [object_to_goal_distance]
8
```

972	9	progress directions = [False] # False indicates that we require
973		\rightarrow the distance to decrease to indicate progress.
974	10	
975	11	<pre>return progress_vars, progress_directions</pre>
976		
977		
978		
979		
980		
981		
982		
983		
984		
985		
986		
987		
988		
909		
001		
002		
003		
994		
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012		
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020		
1027		
1022		
1023		
1024		
IULJ		

1026 A.5 ABLATIONS

1028 No feature library Failed CatchUnderarm code:

1 def progress_function(self):

```
1030
1031
1032
1033
1034
1035
```

```
1036
1037
1038
1039
```

```
# Calculate Euclidean distance between the current object

→ position and the goal position

object_to_goal_distance = torch.norm(self.object_pos -

→ self.goal_pos, dim=-1)

# Calculate the norm of the linear velocity of the object

object_velocity_magnitude = torch.norm(self.object_linvel,

→ dim=-1)

# For stage 1, we want the distance to decrease

# For stage 2, we want the velocity's magnitude to decrease

→ (indicating soft catching)

return [object_to_goal_distance, object_velocity_magnitude],
```

```
\hookrightarrow [False, False]
```

No heuristic discretization Failed SwingCup code:

10/15	_	
1045	1	<pre>def progress_function_no_discretization(self) -> torch.Tensor:</pre>
1040	2	<i># Constants for scales or thresholds</i>
1047	3	POSITION_THRESH = 0.05 # acceptable position error in meters
1048	4	ROTATION_THRESH = 0.1 # acceptable rotation error in radians
1049	5	
1050	6	# Compute distances of hands to their respective cup handles
1051	7	<pre>left_nandle_distance = self.dist(self.left_nand_pos,</pre>
1052	0	\Rightarrow self.cup_tert_namete_pos) right handle distance = self dist(self right hand nos
1053	0	= self cup right handle pos)
1054	9	/ berroup_right_handid_pob/
1055	10	# Compute rotational distance to the goal orientation for the
1055		↔ cup
1050	11	<pre>cup_orientation_error = self.rot_dist(self.object_rot,</pre>
1057		\hookrightarrow self.goal_rot)
1058	12	
1059	13	# Assess position accuracy
1060	14	<pre>left_hand_position_accuracy = (left_handle_distance <</pre>
1061	15	\rightarrow POSITION_IHRESH).IIOdl() right hand position accuracy = (right handle distance <
1062	15	= POSITION THRESH) float()
1063	16	- 1051110A_111000()
1064	17	# Assess rotation accuracy
1065	18	rotation_accuracy = (cup_orientation_error <
1066		\hookrightarrow ROTATION_THRESH).float()
1067	19	
1007	20	<pre># Combine these measures into an overall progress metric (scaled</pre>
1000		$\leftrightarrow to 0-1000)$
1069	21	total_accuracy = leit_nand_position_accuracy *
1070	22	\Rightarrow right_hand_position_accuracy * focation_accuracy
1071	22	\Rightarrow long for discrete bin values
1072	23	, long lot abortoto bin varaoo
1073	24	return progress_bins
1074		
1075		
1076		
1077		

1080 A.6 PROGRESS DISCRETIZATION HEURISTICS

1082 Bi-DexHands

Code logic: for an increasing variable, the min progress value is the value at the start of the episode.
If the variable is greater than zero, the max is tracked as the max of the progress values seen so far. If the variable is less than zero, the max is set to 0. Then progress is rescaled between the min and max.
For a decreasing variable, the complementary logic holds.

1088	1	<pre>def compute_bin_from_progress(self, progress_vars,</pre>			
1089		<pre>→ progress_directions):</pre>			
1090	2	""			
1091	3	Compute the binning from the progress variables and directions			
1092	4	# First, set min/max values if not already set			
1093	6	# Also normalize progress vars in here			
1094	7	for i in range(len(progress_directions)):			
1095	8	<pre>if progress_directions[i]:</pre>			
1006	9	<pre>if 'min' + str(i) not in self.extras:</pre>			
1007	10	<pre>self.extras['min' + str(i)] =</pre>			
1000		\rightarrow torch.min(progress_vars[i])			
1090	11	<pre>if 'max' + str(1) in self.extras: colf outroa[[mau] + str(i)] =</pre>			
1099	12	sell.extras['Mdx' + Str(l)] =			
1100		\Rightarrow colen.max(colen.censor([colen.max(progress_vars[i]),			
1101	13	else:			
1102	14	<pre>self.extras['max' + str(i)] =</pre>			
1103		<pre> → torch.max(progress_vars[i]) </pre>			
1104	15	<pre>if self.extras['max' + str(i)] < 0:</pre>			
1105	16	<pre>progress_vars[i] = torch.clamp((progress_vars[i] -</pre>			
1106		<pre> → self.extras['min' + str(i)]) / </pre>			
1107		→ (-self.extras['min' + str(i)]), min=0, max=1)			
1108	17	erse:			
1100	16	self extras['min' + str(i)]) /			
1110		<pre></pre>			
1110		\rightarrow + str(i)]), min=0, max=1)			
1111	19	else:			
1112	20	<pre>if 'max' + str(i) not in self.extras:</pre>			
1113	21	<pre>self.extras['max' + str(i)] =</pre>			
1114		\rightarrow torch.max(progress_vars[i])			
1115	22	<pre>if 'min' + str(1) in self.extras: colf outroa[[min] + str(i)] =</pre>			
1116	23	self.exclas[min + Stl(1)] =			
1117		$\Rightarrow \text{ self.extras}['\min' + \text{str}(i)]))$			
1118	24	else:			
1119	25	<pre>self.extras['min' + str(i)] =</pre>			
1120		\hookrightarrow torch.min(progress_vars[i])			
1121	26	<pre>if self.extras['max' + str(i)] < 0:</pre>			
1122	27	<pre>progress_vars[i] = torch.clamp((self.extras['max' +</pre>			
1123		<pre> str(1)] - progress_vars[1]), min=0) / (colf ovtroa[]mov[] + str(i)] - colf ovtroa[]min]</pre>			
112/		$\hookrightarrow (\text{Sell}, \text{exclas} [\text{max} + \text{Sell}(1)] = \text{Sell}, \text{exclas} [\text{min}$ $\hookrightarrow + \text{str}(1)])$			
1105	28	else:			
1120	29	progress_vars[i] = torch.clamp((self.extras['max' +			
1120		\rightarrow str(i)] - progress_vars[i]), min=0) /			
1127		<pre> self.extras['max' + str(i)] </pre>			
1128	30	<pre>print("Extras", self.extras)</pre>			
1129	31	<pre># Progress is now associated with increasing values for both</pre>			
1130	22	\rightarrow D105 # So we can generate an overall prograss him by just adding them			
1131	32	π so we can generate an overall progress bin by just adding them \hookrightarrow together, with the appropriate granularity/scaling			
1132	33	binning = torch, zeros (progress vars[0], shape, dtvpe=torch.long.			
1133		→ device=progress_vars[0].device)			

```
1134
                for i in range(len(progress_vars)):
         34
1135
                     binning += ((progress_vars[i] \star (1000 \star (i ==
         35
1136
                     \leftrightarrow (len(progress_vars) - 1)) + 20)).long() % 10000)
1137
                 # Now generate bins from normalized vars
         36
                return binning
         37
1138
1139
       MiniGrid
1140
1141
          1 def discretize_progress(self, obs, max_progress):
1142
                 # Progress is decreasing, max_progress is set as the progress
          2
1143
                \leftrightarrow value at the start of the episode
                # Get progress vars
1144
          3
          4
                progress_vars, _ = self.progress_function()
1145
          5
1146
                # Replace infs and nans in progress_vars
          6
1147
          7
                progress_vars = [0 if math.isnan(var) or math.isinf(var) else var
1148
                 → for var in progress_vars]
1149
          8
          9
                 # Clip by max progress
1150
                if max_progress is None:
         10
1151
         11
                    max_progress = [elem for elem in progress_vars]
1152
         12
                else:
1153
                    progress_vars = [min(var, max_progress[i]) for i, var in
         13

        → enumerate(progress_vars)]

1154
         14
1155
                 # Combine bins
         15
1156
                obs['goal_distance'] = progress_vars[0] +
         16
1157
                 → 100*progress_vars[1]*(progress_vars[0] == 0)
1158
         17
         18
                return obs, max_progress
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
```

1188 A.7 BI-DEXHANDS RESULTS

The results in this section correspond to Fig. 3. The bar chart numbers are reflected in tabular form. All numbers are measured as averages across trials. For ProgressCounts, we also include the standard deviation across 5 trials. For Eureka, we simply report the mean since the standard deviation is not included in the original Eureka paper.

1194			
1195	Task name	ProgressCounts	Eureka
1196	Over	0.93 ± 0.01	0.92
1197	DoorCloseInward	1.00 ± 0.00	1.00
1198	DoorCloseOutward	0.90 ± 0.13	0.96
1199	DoorOpenInward	0.07 ± 0.15	0.00
1200	DoorOpenOutward	0.99 ± 0.01	1.00
1201	Scissors	1.00 ± 0.00	1.00
1202	Swing cup	0.97 ± 0.02	0.66
1203	Switch	0.00 ± 0.00	0.00
1204	Kettle	0.99 ± 0.01	0.89
1205	LiftUnderarm	0.22 ± 0.24	0.70
1206	Pen	0.49 ± 0.06	0.57
1200	BottleCap	0.94 ± 0.03	0.32
1207	CatchAbreast	0.56 ± 0.04	0.50
1208	CatchOver2UnderArm	0.90 ± 0.03	0.90
1209	CatchUnderarm	0.76 ± 0.05	0.67
1210	ReOrientation	0.03 ± 0.00	0.31
1211	GraspAndPlace	0.99 ± 0.01	0.50
1212	BlockStack	0.05 ± 0.05	0.14
1213	PushBlock	0.03 ± 0.03	0.09
1214	TwoCatchUnderarm	0.03 ± 0.02	0.00

Table 8: A comparison of task performance between ProgressCounts and Eureka across the 20
 tasks in Bi-DexHands. Results are averaged across 5 trials for both methods, standard deviation is

1218 reported for ProgressCounts.

1242 A.8 PROGRESS VS PROGRESS DIFFERENCES AS REWARD

While some prior work on reward shaping uses the equivalent of progress differences as reward, we
find that using progress directly as reward leads to better average task success rate than using progress
differences on Bi-DexHands. Therefore, in the main paper we use progress directly as a reward when
comparing to count-based rewards in Section 5.3.

1248			
1249	Task name	ProgressAsReward	ProgressDifferenceAsReward
1250	Average	0.45	0.32
1251		0.00	0.00
1252	Over	0.90	0.88
1253	DoorCloseInward	1.00	0.96
1254	DoorCloseOutward	1.00	0.00
1234	DoorOpenInward	0.00	0.31
1255	DoorOpenOutward	0.31	0.00
1256	Scissors	1.00	1.00
1257	SwingCup	0.99	0.94
1258	Switch	0.00	0.00
1259	Kettle	0.00	0.04
1260	LiftUnderarm	0.08	0.35
1261	Pen	0.22	0.00
1262	BottleCap	0.04	0.00
1263	CatchAbreast	0.49	0.01
1064	CatchOver2UnderArm	0.94	0.00
1204	CatchUnderarm	0.88	0.82
1265	ReOrientation	0.06	0.03
1266	GraspAndPlace	0.98	1.00
1267	BlockStack	0.00	0.00
1268	PushBlock	0.02	0.00
1269	TwoCatchUnderarm	0.01	0.00
1270			

Table 9: Using progress directly as reward leads to better average task success rate than using progress differences on Bi-DexHands.



Figure 5: Training curves for ProgressCounts on 8 hard-exploration MiniGrid tasks.

A.10 EXPERIMENTAL DETAILS AND HYPERPARAMETERS

Environments We train our policies on two environments: Bi-DexHands (Chen et al., 2022), and MiniGrid (Chevalier-Boisvert et al., 2024).

Code For the Bi-DexHands benchmark, we build upon the codebase from Eureka (Ma et al., 2023): https://github.com/eureka-research/Eureka. The repo uses the RLGames implementation of PPO for training (Makoviichuk & Makoviychuk, 2021). For the MiniGrid benchmark, we build upon the codebase from NovelD (Zhang et al., 2021): https://github.com/tianjunz/NovelD. Our experimental code is is avail-able at the following anonymous link: https://drive.google.com/drive/folders/ 1G88Je0K4BuexWhE8ZLgoM0WM6vHcBrvt?usp=sharing.

Hyperparameters For all Bi-DexHands tasks, we scale extrinsic rewards by 0.05, and normalize intrinsic rewards to a mean of 0.001. Elsewhere, we use the default hyperparameters associated with Eureka, which are the default parameters from the original Bi-DexHands benchmark. Progress functions are discretized into 1020 bins for count-based exploration (for tasks with two subtasks, the first subtask is discretized to 20 bins, and the second subtask to 1000 bins). For all MiniGrid tasks, we set an intrinsic reward coefficient of 0.5, and leave all other hyperparameters at default values. Progress functions are discretized into 50 bins for count-based exploration (for tasks with two subtasks, each subtask is discretized to 25 bins).

Compute All experiments were run on a machine with 8 NVidia Tesla V100 GPUs across 3 weeks.