
Exploration via Elliptical Episodic Bonuses

Mikael Henaff
Meta AI Research
mikaelhenaff@meta.com

Roberta Raileanu
Meta AI Research
raileanu@meta.com

Minqi Jiang
University College London
Meta AI Research
meta@fb.com

Tim Rocktäschel
University College London
t.rocktaschel@cs.ucl.ac.uk

Abstract

In recent years, a number of reinforcement learning (RL) methods have been proposed to explore complex environments which differ across episodes. In this work, we show that the effectiveness of these methods critically relies on a count-based episodic term in their exploration bonus. As a result, despite their success in relatively simple, noise-free settings, these methods fall short in more realistic scenarios where the state space is vast and prone to noise. To address this limitation, we introduce **Exploration via Elliptical Episodic Bonuses (E3B)**, a new method which extends count-based episodic bonuses to continuous state spaces and encourages an agent to explore states that are diverse under a learned embedding within each episode. The embedding is learned using an inverse dynamics model in order to capture controllable aspects of the environment. Our method sets a new state-of-the-art across 16 challenging tasks from the MiniHack suite, without requiring task-specific inductive biases. E3B also matches existing methods on sparse reward, pixel-based Vizdoom environments, and outperforms existing methods in reward-free exploration on Habitat, demonstrating that it can scale to high-dimensional pixel-based observations and realistic environments.

1 Introduction

Exploration in environments with sparse rewards is a fundamental challenge in reinforcement learning (RL). In the tabular setting, provably optimal algorithms have existed since the early 2000s [36, 10]. More recently, exploration has been studied in the context of deep RL, and a number of empirically successful methods have been proposed, such as pseudocounts [9], intrinsic curiosity modules (ICM) [52], and random network distillation (RND) [13]. These methods rely on intrinsically generated exploration bonuses that reward the agent for visiting states that are novel according to some measure. Different measures of novelty have been proposed, such as the likelihood of a state under a learned density model, the error of a forward dynamics model, or the loss on a random prediction task. These approaches have proven effective on hard exploration problems, as exemplified by the Atari games Montezuma’s Revenge and PitFall [22].

The approaches above are, however, designed for *singleton* RL tasks, where the agent is spawned in the same environment in every episode. Recently, several studies have drawn attention to the fact that RL agents exhibit poor generalization across environments, and that even minor changes to the environment can lead to substantial degradation in performance [75, 35, 74, 18, 38]. This has motivated the creation of benchmarks in the Contextual Markov Decision Process (CMDP) framework, where different episodes correspond to different environments that nevertheless share certain characteristics. Examples of CMDPs include procedurally generated (PCG) environments

[15, 58, 41, 34, 17, 7, 29, 53] or embodied AI tasks where the agent must generalize its behavior to unseen physical spaces at test time [59, 61, 28, 72]. A number of methods have been proposed which have shown promising performance in PCG environments with sparse rewards, such as RIDE [56], AGAC [25] and NovelD [76]. These methods propose different intrinsic reward functions, such as the change in representation in a latent space, the divergence between the predictions of a policy and an adversary, or the difference between random network prediction errors at two consecutive states. Although not presented as a central algorithmic feature, these methods also include a count-based bonus which is computed at the episode level.

In this work, we take a closer look at exploration in CMDPs, where each episode corresponds to a different environment context. We first show that, surprisingly, the count-based episodic bonus that is often included as a heuristic is in fact essential for good performance, and current methods fail if it is omitted. Furthermore, due to this dependence on a count-based term, existing methods fail on more complex tasks with irrelevant features or dynamic entities, where each observation is rarely seen more than once. We find that performance can be improved by counting certain features extracted from the observations, rather than the observations themselves. However, different features are useful for different tasks, making it difficult to design a feature extractor that performs well across all tasks.

To address this fundamental limitation, we propose a new method, E3B, which uses an elliptical bonus [6, 20, 42] at the episode level that can be seen as a natural generalization of a count-based episodic bonus to continuous state spaces, and that is paired with a self-supervised feature learning method using an inverse dynamics model. Our algorithm is simple to implement, scalable to large or infinite state spaces, and achieves state-of-the-art performance across 16 challenging tasks from the MiniHack suite [58], without the need for task-specific prior knowledge. It also matches existing methods on hard exploration tasks from the VizDoom environment [37], and significantly outperforms existing methods in reward-free exploration on the Habitat embodied AI environment [59, 69]. This demonstrates that E3B scales to rich, high dimensional pixel-based observations and real-world scenarios. Our code is available at <https://github.com/facebookresearch/e3b>.

2 Background

2.1 Contextual MDPs

We consider a Contextual Markov Decision Process ¹ (CMDP [30]) given by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{C}, P, r, \mu_C, \mu_S)$ where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{C} is the context space, P is the transition function, r is the reward function, μ_C is the distribution over contexts and μ_S is the conditional initial state distribution. At each episode, we sample a context $c \sim \mu_C$, an initial state $s_0 \sim \mu(\cdot|c)$, and subsequent states in the episode are sampled according to $s_{t+1} \sim P(\cdot|s_t, a_t, c)$. Let d_π^c denote the distribution over states induced by following policy π in context c . The goal is to learn a policy π which maximizes the expected return over all contexts, i.e. $R = \mathbb{E}_{c \sim \mu_C, s \sim d_\pi^c, a \sim \pi(s)}[r(s, a)]$. Examples of CMDPs include procedurally-generated environments, such as ProcGen [17], MiniGrid [15], NetHack [41], or MiniHack [58], where each context c corresponds to the random seed used to generate the environment; in this case, the number of contexts $|\mathcal{C}|$ is effectively infinite. Other examples include embodied AI environments [59, 69, 28, 61, 72], where the agent is placed in different simulated houses and must navigate to a location or find an object. In this setting, each context $c \in \mathcal{C}$ represents a house identifier and the number of houses $|\mathcal{C}|$ is typically between 20 and 1000. For an in-depth review of the literature on CMDPs and generalization in RL, see [39].

2.2 Exploration Bonuses

If the environment rewards are sparse, learning a policy using simple ϵ -greedy exploration may require intractably many samples. We therefore consider methods that augment the external reward function r with an intrinsic reward bonus b . A number of intrinsic bonuses that encourage exploration in singleton (non-contextual) MDPs have been proposed, including pseudocounts [9], intrinsic curiosity modules (ICM) [52] and random network distillation (RND) error [13]. At a high level, these methods

¹Technically, some of the environments we consider are Contextual Partially Observed MDPs (CPOMDPs), but we follow the convention in [39] and adopt the CMDP framework for simplicity. For CPOMDPs, we use recurrent networks or frame stacking to convert to CMDPs, as done in prior work [46].

define an intrinsic reward bonus that is high if the current state is different from the previous states visited by the agent, and low if it is similar (according to some measure).

Method	Exploration bonus
RIDE [56]	$\ \phi(s_{t+1}) - \phi(s_t)\ _2 \cdot 1/\sqrt{N_e(s_{t+1})}$
NovelD [76]	$\left[b_{\text{RND}}(s_{t+1}) - \alpha \cdot b_{\text{RND}}(s_t) \right]_+ \cdot \mathbb{I}[N_e(s_{t+1}) = 1]$
AGAC [25]	$D_{\text{KL}}(\pi(\cdot s_t) \ \pi_{\text{adv}}(\cdot s_t)) + \beta \cdot 1/\sqrt{N_e(s_{t+1})}$

Table 1: Summary of recent exploration methods for procedurally-generated environments. Each exploration bonus has a count-based episodic term, marked in blue.

More recently, several methods have been proposed for and evaluated on procedurally-generated MDPs. All use different exploration bonuses, which are summarized in Table 1. RIDE [56] defines a bonus based on the distance between the embeddings of two consecutive observations, NovelD [76] uses a bonus based on the difference between two consecutive RND bonuses, and AGAC [25] uses the KL divergence between the predictions of the agent’s policy and those of an adversarial policy trained to mimic it (see the original works for details). In addition, all three methods have a term in their bonus (marked in blue) which depends on $N_e(s_{t+1})$ – the number of times s_{t+1} has been encountered during the *current episode*. Although presented as a heuristic, below we will show that without this count-based episodic term, all three methods fail to learn. This in turn limits their effectiveness in more complex, dynamic, and noisy environments.

3 Importance and Limitations of Count-Based Episodic Bonuses

We now discuss in more detail the importance and limitations of the count-based episodic bonuses used in RIDE, AGAC and NovelD, which depend on $N_e(s_t)$. Figure 1 shows results for the three methods with and without their respective count-based episodic terms, on one of the MiniGrid environments used in prior work. When the count-based terms are removed, all three methods fail to learn. Similar trends apply for other MiniGrid environments (see Appendix D.1). This shows that the episodic bonus is in fact essential for good performance.

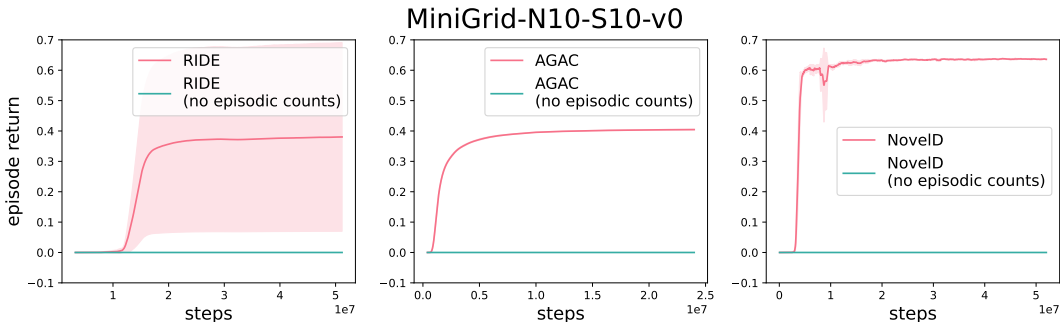


Figure 1: Results for RIDE, AGAC and NovelD with and without the count-based episodic bonus, over 5 random seeds (solid line indicates the mean, shaded region indicates one standard deviation). All three algorithms fail if the count-based episodic bonus is removed.

However, the count-based episodic bonus suffers from a fundamental limitation, which is similar to that faced by count-based approaches in general: if each state is unique, then $N_e(s_t)$ will always be 1 and the episodic bonus is no longer meaningful. This is the case for many real-world applications. For example, a household robot’s state as recorded by its camera might include moving trees outside the window, clocks showing the time or images on a television screen which are not relevant for its tasks, but nevertheless make each state unique.

Previous works [56, 25, 76] have used the MiniGrid test suite [15] for evaluation, where observations are less noisy and do not typically contain irrelevant information. Thus, methods relying on episodic counts have been effective in these scenarios. However, in more complex environments such as

MiniHack [58] or with high-dimensional pixel-based observations, episodic count-based approaches can cease to be viable.

A possible alternative could be to design a function to extract relevant features from each state and feed them to the count-based episodic bonus. Specifically, instead of defining a bonus using $N_e(s_t)$, we could define the bonus using $N_e(\phi(s_t))$, where ϕ is a hand-designed feature extractor. For example, in the paper introducing the MiniHack suite [58], the RIDE implementation uses $\phi(s_t) = (x_t, y_t)$, where (x_t, y_t) is the spatial location of the agent at time t . However, this approach relies heavily on task-specific knowledge.

Figure 2 shows results on two tasks from the MiniHack suite for three NovelD variants (we decided to focus our study on variants of NovelD, since it was previously shown to outperform competing methods on MiniGrid [76]). The first, NOVELD, denotes the standard formulation which uses the bonus $\mathbb{I}[N_e(s_t) = 1]$. The second, NOVELD-POSITION, uses the positional feature encoding described above, i.e. the episodic bonus is defined as $\mathbb{I}[N_e(\phi(s_t)) = 1]$ with $\phi(s_t) = (x_t, y_t)$. The third, NOVELD-MESSAGE, uses a feature encoding where $\phi(s_t)$ extracts the message portion of the state s_t similarly to [47] (both encodings are explained in more detail in Section 5). In contrast to the MiniGrid environments, here standard NOVELD fails completely due to the presence of a time counter feature in the MiniHack observations, which makes each observation in the episode unique. Using the positional encoding enables NOVELD-POSITION to solve the MultiRoom task, but this method fails on the Freeze task. On the other hand, using the message encoding enables NOVELD-MESSAGE to succeed on the Freeze task, but it fails on the MultiRoom one. This illustrates that different feature extractors are effective on different tasks, and that designing one which is broadly effective is challenging. Therefore, robust new methods which do not require task-specific engineering are needed.

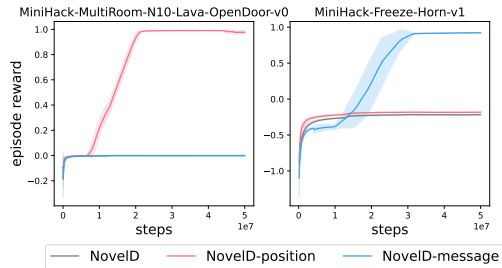


Figure 2: Performance of NovelD with different feature extractors: entire observation, agent location, or environment message. Results are averaged over 5 seeds and the shaded region represents one standard deviation.

4 Elliptical Episodic Bonuses

In this section we describe **Exploration via Elliptical Episodic Bonuses, (E3B)**, our algorithm for exploration in contextual MDPs. It is designed to address the shortcomings of count-based episodic bonuses described above, with two aims in mind. First, we would like an episodic bonus that can be used with *continuous* state representations, unlike the count-based bonus which requires discrete states. Second, we would like a representation learning method that only captures information about the environment that is relevant for the task at hand. The first requirement is met by using an elliptical bonus [6, 20, 42], which provides a continuous analog to the count-based bonus, while the second requirement is met by using a representation learned with an inverse dynamics model [52, 56].

A summary of the method is shown in Figure 3. We define an intrinsic reward based on the position of the current state’s embedding with respect to an ellipse fit on the embeddings of previous states encountered within the same episode. This bonus is then combined with the environment reward and used to update the agent’s policy. The next two sections describe the elliptical bonus and embedding method in detail.

4.1 Elliptical Episodic Bonus

Given a feature encoding ϕ , at each time step t in the episode the elliptical bonus b is defined as follows:

$$b(s_t) = \phi(s_t)^\top C_{t-1}^{-1} \phi(s_t), \quad C_{t-1} = \sum_{i=1}^{t-1} \phi(s_i) \phi(s_i)^\top + \lambda I \quad (1)$$

Here λI is a regularization term to ensure that the matrix C_{t-1} is non-singular, where λ is a scalar coefficient and I is the identity matrix. The reward optimized by the algorithm is then defined as

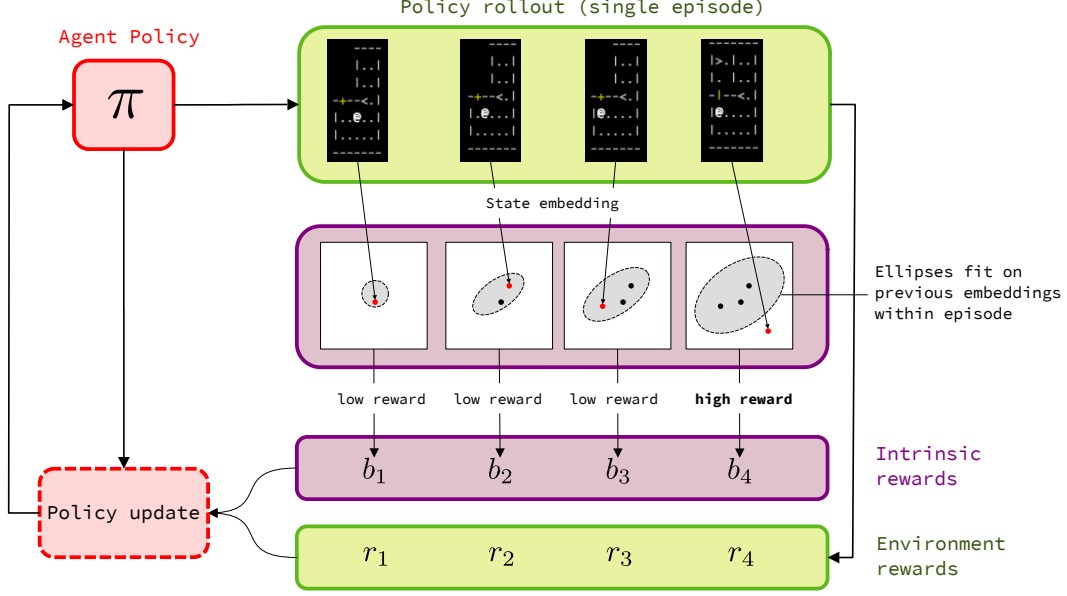


Figure 3: Overview of E3B. At each step in an episode, an ellipse is fit on the embeddings of previous states encountered within the episode. Intrinsic rewards are computed based on the position of the current state’s embedding with respect to the ellipse: the further outside the ellipse, the higher the reward. These are combined with environment rewards, the policy is updated, and the process repeated. The state embeddings are learned using an inverse dynamics model.

$\bar{r}(s_t, a_t) = r(s_t, a_t) + \beta \cdot b(s_t)$, where $r(s_t, a_t)$ is the extrinsic reward provided by the environment and β is a scalar term balancing the tradeoff between exploration and exploitation.

Intuition. One perspective is that the elliptical bonus is a natural generalization of a count-based episodic bonus [2]. To see this, observe that if the problem is tabular and ϕ is a one-hot encoding of the state, then C_{t-1} will be a diagonal matrix whose entries contain the counts corresponding to each state encountered in the episode. Its inverse C_{t-1}^{-1} will also be a diagonal matrix whose entries are inverse state visitation counts, and the bilinear form $\phi(s_t)^\top C_{t-1}^{-1} \phi(s_t)$ reads off the entry corresponding to the current state s_t , yielding a bonus of $1/N_e(s_t)$.

For a more general geometric interpretation, if $\phi(s_0), \dots, \phi(s_{t-1})$ are roughly centered at zero, then C_{t-1} can be viewed as their unnormalized covariance matrix. Now consider the eigendecomposition $C_{t-1} = U^\top \Lambda U$, where Λ is the diagonal matrix whose entries are the eigenvalues $\lambda_1, \dots, \lambda_n$ (these are real since C_{t-1} is symmetric). Letting $z = U\phi(s_t) = (z_1, \dots, z_n)$ be the set of coordinates of $\phi(s_t)$ in the eigenspace of C_{t-1} , we can rewrite the elliptical bonus as:

$$b(s_t) = z^\top \Lambda^{-1} z = \sum_{i=1}^n \frac{z_i^2}{\lambda_i}$$

The bonus increases the more $\phi(s_t)$ is aligned with the eigenvectors corresponding to smaller eigenvalues of C_{t-1} (directions of low data density), and decreases the more it is aligned with eigenvectors corresponding to larger eigenvalues (directions of high data density). An illustration for $n = 2$ is shown in Figure 10 of Appendix B. In our experiments, n is typically between 256 and 1024.

Efficient Computation. The matrix C_{t-1} needs to be inverted at each step t in the episode, an operation which is cubic in the dimension of ϕ and may thus be expensive. To address this, we use the Sherman-Morrison matrix identity [62] to perform fast rank-1 updates in quadratic time:

$$C_t^{-1} = \begin{cases} \frac{1}{\lambda} I & \text{if } t = 0 \\ C_{t-1}^{-1} - \frac{C_{t-1}^{-1} \phi(s_t) \phi(s_t)^\top C_{t-1}^{-1}}{1 + \phi(s_t)^\top C_{t-1}^{-1} \phi(s_t)} & \text{if } t \geq 1 \end{cases}$$

This results in an approximately $3\times$ speedup over naïve matrix inversion (details in Appendix D.2).

4.2 Learned Feature Encoder

Any feature learning method could in principle be used to learn ϕ . Here we use the inverse dynamics model approach proposed in [52], which trains a model g along with ϕ to map each pair of consecutive embeddings $\phi(s_t), \phi(s_{t+1})$ to a distribution over actions a_t linking them. In our setup, ϕ is separate from the policy network. The g model is trained jointly with ϕ using the following per-sample loss:

$$\ell(s_t, a_t, s_{t+1}; \phi, g) = -\log p(a_t | g(\phi(s_t), \phi(s_{t+1})))$$

The motivation is that the mapping ϕ will discard information about the environment which is not useful for predicting the agent’s actions. Previous work [52] has shown that this can make learning more robust to random noise or other parts of the state which are not controllable by the agent.

4.3 Full Algorithm

Putting all of these together, the full algorithm is given below.

Algorithm 1 Exploration via Episodic Elliptical Bonuses (E3B)

Initialize policy π , feature encoder ϕ and inverse dynamics model f .

while not converged **do**

 Sample context $c \sim \mu_C$ and initial state $s_0 \sim \mu_S(\cdot | c)$

 Initialize inverse covariance matrix: $C_0^{-1} = \frac{1}{\lambda}I$

for $t = 0, \dots, T$ **do**

$a_t \sim \pi(\cdot | s_t)$ // Sample action

$s_{t+1}, r_{t+1} \sim P(\cdot | s_t, a_t)$ // Step through environment

$b_{t+1} = \phi(s_{t+1})^\top C_t^{-1} \phi(s_{t+1})$ // Compute bonus

$u = C_t^{-1} \phi(s_{t+1})$

$C_{t+1}^{-1} = C_t^{-1} - \frac{1}{1+b_{t+1}} uu^\top$ // Update inverse covariance matrix

$\bar{r}_{t+1} = r_{t+1} + \beta b_{t+1}$

end for

 Perform policy gradient update on π using rewards $\bar{r}_1, \dots, \bar{r}_T$.

 Update ϕ and g using $\{(s_t, a_t, s_{t+1})\}_{t=0}^{T-1}$ to minimize the loss:

$$\ell = -\log(p(a_t | f(\phi(s_t), \phi(s_{t+1}))))$$

end while

5 Experiments

5.1 MiniHack Suite

In order to probe the capabilities of existing methods and evaluate E3B, we seek CMDP environments which exhibit challenges associated with realistic scenarios, such as sparse rewards, noisy or irrelevant features, and large state spaces. For our first experimental testbed, we opted for the procedurally generated tasks from the MiniHack suite [58], which is itself based on the NetHack Learning Environment [41]. NetHack is a notoriously challenging roguelike video game where the agent must navigate through procedurally generated dungeons to recover a magical amulet. The MiniHack tasks contain numerous challenges such as finding and using magical objects, navigating through levels while avoiding lava, and fighting monsters. Furthermore, rewards are sparse and as detailed below, the state representation contains a large amount of information, only some of which is relevant for a given task.

For all our experiments we use the Torchbeast [40] implementation of IMPALA [23] as our base RL algorithm. For certain skill-based tasks, we restricted the action space to the necessary actions for solving the task at hand, since we found that none of the methods were able to make progress with

the full action space (see Appendix C.1.4). See Appendix C.1.3 for environment details and C.1 for other experiment details.

5.1.1 Modalities for Episodic Bonus

The MiniHack environments provide an observation at each step that includes three different modalities: i) a symbolic image, which indicates the location of the agent, monsters or other entities, objects and different types of terrain (such as walls, lava, water); ii) a statistics vector which indicates the agent’s current (x, y) location, hit points, time step t , and other features such as strength, dexterity and constitution, and iii) a textual message which gives different types of feedback about the environment ("you see a key of master thievery", "the wall is solid rock"). These are illustrated in Figure 4.

We now draw attention to some important subtleties regarding how existing methods implement count-based episodic bonuses—as we will see, these can have a large impact on performance. The works of [56, 25, 76], which use the MiniGrid suite [15] for evaluation, use a hash table whose keys are the full observations (in this case symbolic images) encountered by the agent. The occurrences of these observations are then counted to compute the episodic bonus. The baselines run in the MiniHack tasks [58], on the other hand, use a hash table whose keys are the (x, y) positions of the agent, which are extracted from the observation and are then similarly counted to compute the episodic bonus.

In order to understand the effect of these different input modalities on count-based episodic bonuses, we consider three variants of the NovelD algorithm, in addition to the original formulation. These methods all have a count-based episodic bonus of the form $\mathbb{I}[N_e(\phi(s_t)) = 1]$, with different choices of ϕ detailed below.

NOVELD: in the original version of the algorithm, ϕ is the identity. Here the input to the count-based episodic bonus is the full state, namely the concatenation of the symbol image, the message and the stats vector. This makes no assumptions about which part of the state is most useful.

NOVELD-POSITION: in this variant, ϕ extracts the (x, y) position of the agent from the statistics vector, which reflects a strong inductive bias that the task is navigation-based.

NOVELD-IMAGE: in this variant, ϕ extracts the symbol image only. This reflects an inductive bias that the message and the stats vector are not useful since they are discarded.

NOVELD-MESSAGE: in this variant, ϕ extracts the message only. This reflects an inductive bias that the messages are important but the stats and symbolic images are not.

For E3B, we feed the full state to the algorithm and do not make any assumptions about which part is useful for the task at hand. Despite the lack of task-specific inductive biases, our method is still able to extract the relevant features for each task, thus outperforming the other exploration approaches (or matching their performance when prior knowledge of the task is used).

5.1.2 Results on MiniHack

Aggregate results for IMPALA, RND, RIDE, ICM, NovelD (with the three variants described above) and E3B over 16 sparse reward tasks from the MiniHack suite are shown in Figure 5. Out of 16 environments, 8 are based on the MiniGrid suite, but use the MiniHack interface and observation space (environment details are included in Appendix C.1.3). We used the performance metrics and bootstrapping protocol from [4] to compute confidence intervals, which are more informative than simple point estimates. We see that over all tasks (top row), E3B outperforms all other methods by a significant margin across all three performance metrics.

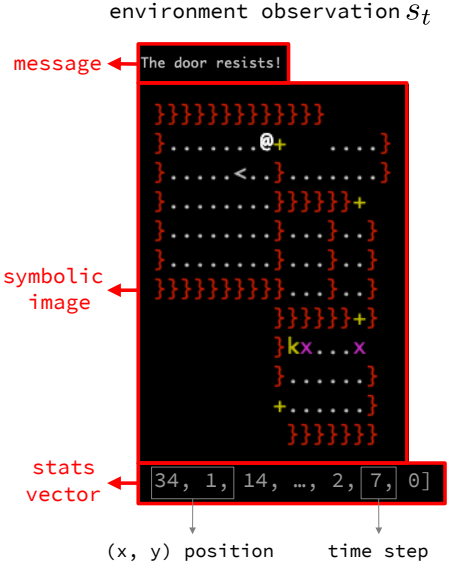


Figure 4: Observation for MiniHack

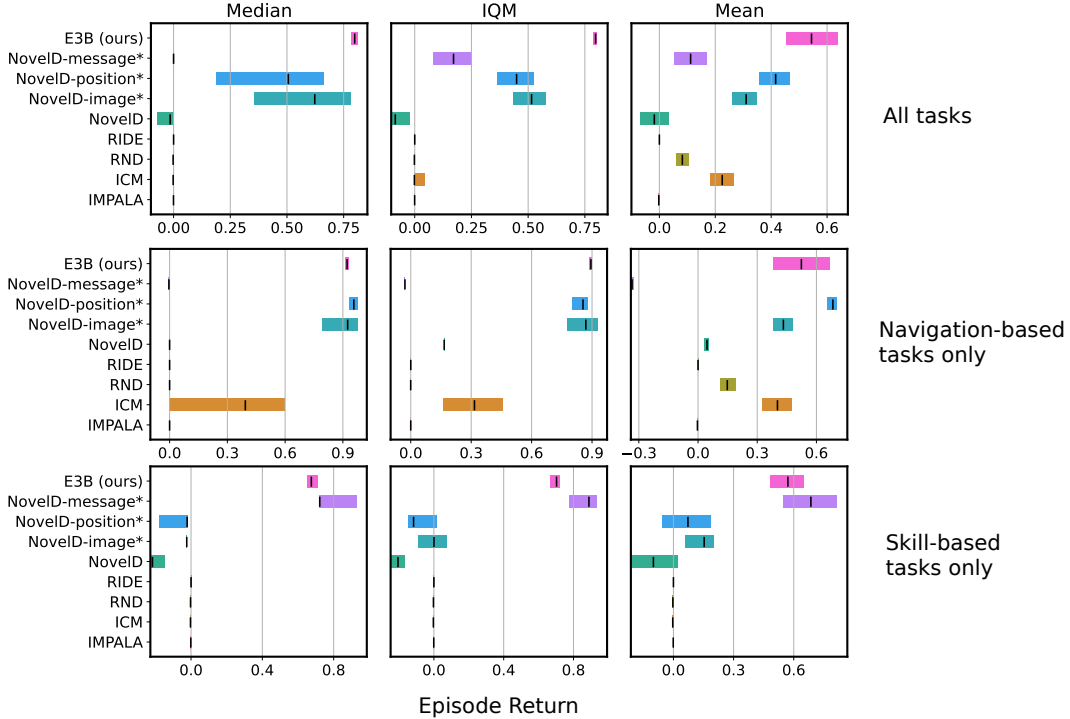


Figure 5: Aggregate results over 16 tasks from the MiniHack environment. Bars represent 95% confidence intervals computed using stratified bootstrapping with 5 random seeds. Methods marked with * use task-specific prior knowledge.

Standard NOVELD performs poorly due to the fact that the MiniHack observations contain a time counter in the statistics vector (see Figure 4), which makes each observation in the episode unique and hence renders the count-based episodic bonus meaningless. The three NovelD variants which use different features extracted from the state for the count-based bonus perform better. Figure 5 shows aggregate performance across a subset of 9 navigation-based tasks (middle row) and 7 skill-based tasks (bottom row) (see Appendix C.1.3 for the task breakdown). On the navigation-based tasks, NOVELD-POSITION has excellent performance, since visiting a large number of different (x, y) locations is closely aligned with the true task reward. However, NOVELD-POSITION fails on all the skill-based tasks, resulting in poor performance overall. We observe an opposite trend for the NOVELD-MESSAGE variant, which performs very well on the skill-based tasks, but poorly on the navigation-based ones.

This highlights that although certain inductive biases can help for certain tasks, it is difficult to find one which performs well across all of them. Our method, E3B, performs well on both the navigation-based tasks and the skill-based tasks, resulting in superior performance overall. It is worth noting that E3B does not require any task-specific engineering: the exact same algorithm is run for all tasks, and it uses the unprocessed states. Results for individual tasks, along with more analysis and discussion, can be found in Appendix D.4.

5.1.3 Ablation Experiments

We next report results for ablation experiments measuring the effects of different algorithmic components of E3B, namely the feature encoding ϕ and the episodic nature of the bonus. Results are shown in Figure 6. E3B (random enc.) indicates E3B where ϕ is a randomly initialized network which is kept fixed throughout training. E3B (policy enc.) indicates E3B where the weights of ϕ are tied to those of the policy network (with the last layer producing action probabilities removed). E3B (non-episodic) indicates E3B where the

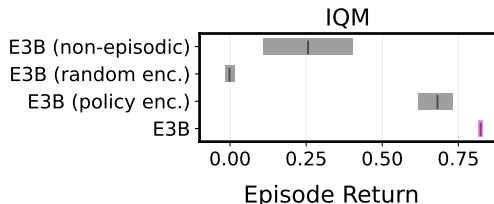


Figure 6: MiniHack ablations, bars represent 95% confidence intervals using stratified bootstrapping.

elliptical bonus is computed using observations across *all* timesteps in the lifetime of the agent, and not just the current episode. All three variants perform significantly worse than E3B, which uses an inverse dynamics model encoding for ϕ and an episodic bonus. This highlights that both the inverse dynamics model and episodic bonus are important for success.

5.2 Pixel-Based VizDoom

As our second evaluation testbed, we used the sparse reward, pixel-based VizDoom [37] environments used in prior work [52, 56]. Although these are singleton MDPs, they still constitute challenging exploration problems and probe whether our method scales to continuous high-dimensional pixel-based observations. Results comparing E3B to RIDE, ICM and IMPALA on three versions of the task are shown in Figure 7 (hyperparameters can be found in Appendix C.2). IMPALA succeeds on the dense reward task but fails on the two sparse reward ones. E3B is able to solve both versions of the sparse reward task, similar to RIDE and ICM.

We emphasize that these are *singleton* MDPs, where the environment does not change from one episode to the next. Therefore, it is unsurprising that ICM, which was designed for singleton MDPs, succeeds in this task. RIDE is also able to solve the task, consistent with results from prior work [56]. The fact that E3B also succeeds provides evidence of its robustness and its applicability to settings with pixel-based observations.

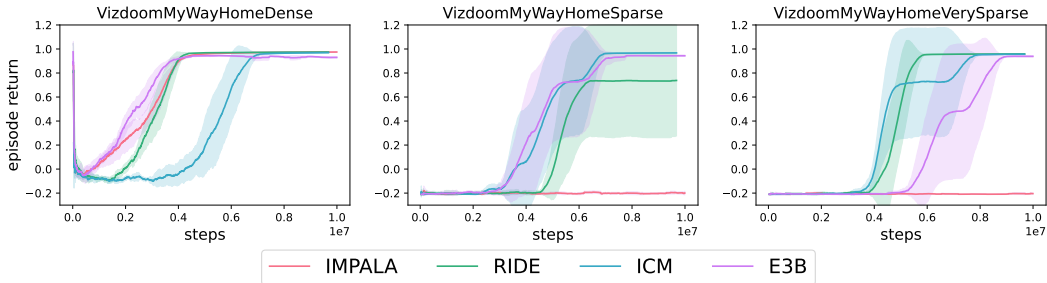


Figure 7: Results on pixel-based Vizdoom tasks with dense and sparse rewards. Results are averaged over 5 random seeds, shaded region indicates one standard deviation.

5.3 Reward-free Exploration on Habitat

As our third experimental setting, we investigate reward-free exploration in Habitat [59, 69]. Habitat is a platform for embodied AI research which provides an interface for agents to navigate and act in photorealistic simulations of real indoor environments. At each episode during training, the agent is initialized in a different environment, and it is tested on a set of held-out environments not used during training. These experiments are designed to evaluate exploration of realistic CMDPs with visually rich observations.

We use the HM3D dataset [57], which contains high-quality renditions of 1000 different indoor spaces. As our base RL algorithm we use DD-PPO [71] and train ICM, RND, NovelD and E3B agents using the intrinsic reward alone. We then evaluate each agent (as well as two random agents) on unseen test environments by measuring how much of each environment has been revealed by the agent’s line of sight over the course of the episode. Full details on the experimental setup can be found in Appendix C.3.

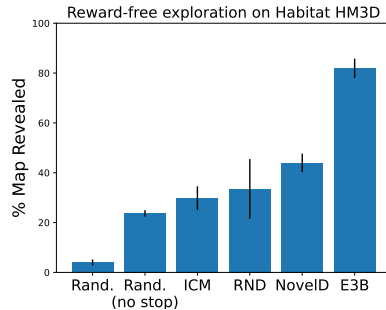


Figure 8: Reward-free exploration on Habitat. Error bars represent std. deviations over 3 seeds.

Quantitative results are shown in Figure 8. The E3B agent reveals significantly more of the test maps than any of the other agents. Trajectories for E3B, ICM, RND and NovelD on one of the test maps are shown in Figure 9, which illustrate how E3B explores a large portion of the space whereas the

others do not. These results provide evidence for E3B’s scalability to high-dimensional pixel-based observations, and reinforce its broad applicability.

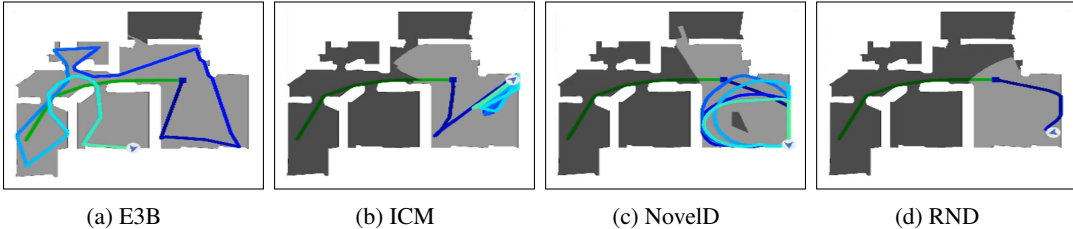


Figure 9: Trajectories of policies trained with different exploration algorithms, on a Habitat environment unseen during training. E3B reveals a larger portion of the map than other methods.

6 Related Work

Exploration in RL. Exploration remains a long-standing problem in RL. Common approaches include ϵ -greedy [68], count-based exploration [66, 8, 48, 45, 70, 43], curiosity-based exploration [60, 64, 65, 11], and other types of intrinsic motivation [49, 50, 63, 1, 12]. These methods were largely designed for singleton MDPs, where the environment remains the same across episodes. As a result, they measure novelty over the agent’s lifetime in a static environment rather than on an episodic basis under a distribution of diverse contexts, as necessary in CMDPs. Other intrinsic motivation methods have recently been developed for exploration in CMDPs [56, 73, 25, 76]. As our experiments show, they critically rely on episodic count-based bonuses, which E3B generalizes.

Another class of methods automatically generate curricula over variations of the CMDP to encourage efficient learning, effectively performing a form of curiosity-driven exploration in the context space. These include goal-conditioned [27, 26, 24, 55, 19, 14, 22] and goal-free variants [67, 54, 33, 21]. Unlike our method, these methods assume the ability to actively configure the environment context. In principle, E3B can be combined with these approaches, whereby E3B explores at an episodic level, while the automatic curriculum method explores at a context level, making the two complementary.

Elliptical Bonuses. The use of elliptical bonuses has a long history in the contextual bandit literature [6, 20, 42], which corresponds to the RL setting with a single time step. More recently, several works have begun to explore the use of elliptical bonuses in the context of multi-step RL. The PC-PG algorithm [2] considers MDPs with linear dynamics and uses an elliptical bonus based on a policy cover to explore in a provably efficient manner. The ACB algorithm [5] also uses an elliptical bonus, approximated using linear regressors trained on random noise, while FLAMBE [3] uses an elliptical bonus inside a learned dynamics model. All of these algorithms operate on singleton MDPs, whereas ours is designed for contextual MDPs and constructs elliptical bonuses at the episode level rather than across episodes. We also use different feature learning methods, namely inverse dynamics models, instead of the policy encoders, random networks or kernel methods used in the aforementioned works.

7 Conclusion

In this work, we identified a fundamental limitation of existing methods for exploration in CMDPs: their performance relies heavily on an episodic count-based term, which is not meaningful when each state is unique. This is a common scenario in realistic applications, and it is difficult to alleviate through feature engineering. To remedy this limitation, we introduce a new method, E3B, which extends episodic count-based methods to continuous state spaces using an elliptical episodic bonus, as well as an inverse dynamics model to automatically extract useful features from states. E3B achieves a new state-of-the-art on a wide range of complex tasks from the MiniHack suite, without the need for feature engineering. Our approach also scales to high-dimensional pixel-based environments, demonstrated by the fact that it matches top exploration methods on VizDoom and outperforms them in reward-free exploration on Habitat. Future research directions include experimenting with more advanced feature learning methods, and investigating ways to integrate within-episode and across-episode novelty bonuses.

References

- [1] Joshua Achiam and Shankar Sastry. Surprise-based intrinsic motivation for deep reinforcement learning. *arXiv preprint arXiv:1703.01732*, 2017.
- [2] Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. Pc-pg: Policy cover directed exploration for provable policy gradient learning. *Advances in neural information processing systems*, (33).
- [3] Alekh Agarwal, Sham M. Kakade, Akshay Krishnamurthy, and Wen Sun. FLAMBE: structural complexity and representation learning of low rank mdps. *CoRR*, abs/2006.10814, 2020.
- [4] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Information Processing Systems*, 2021.
- [5] Jordan T. Ash, Cyril Zhang, Surbhi Goel, Akshay Krishnamurthy, and Sham M. Kakade. Anti-concentrated confidence bonuses for scalable exploration. *CoRR*, abs/2110.11202, 2021.
- [6] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3:397–422, March 2002.
- [7] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Victor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016.
- [8] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- [9] Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, page 1479–1487, Red Hook, NY, USA, 2016. Curran Associates Inc.
- [10] Ronen I. Brafman and Moshe Tennenholtz. R-MAX - A general polynomial time algorithm for near-optimal reinforcement learning. *J. Mach. Learn. Res.*, 3:213–231, 2002.
- [11] Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A. Efros. Large-scale study of curiosity-driven learning. In *ICLR*, 2019.
- [12] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- [13] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.
- [14] Andres Campero, Roberta Raileanu, Heinrich Küttler, Joshua B Tenenbaum, Tim Rocktäschel, and Edward Grefenstette. Learning with amigo: Adversarially motivated intrinsic goals. *arXiv preprint arXiv:2006.12122*, 2020.
- [15] Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld environment for openai gym. <https://github.com/maximecb/gym-minigrid>, 2018.
- [16] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [17] Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2048–2056. PMLR, 13–18 Jul 2020.
- [18] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1282–1289. PMLR, 09–15 Jun 2019.
- [19] Cédric Colas, Tristan Karch, Olivier Sigaud, and Pierre-Yves Oudeyer. Intrinsically motivated goal-conditioned reinforcement learning: a short survey. *ArXiv*, abs/2012.09830, 2020.

- [20] Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic linear optimization under bandit feedback. In *COLT*, 2008.
- [21] Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. *Advances in Neural Information Processing Systems*, 33:13049–13061, 2020.
- [22] Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Go-explore: a new approach for hard-exploration problems. *arXiv preprint arXiv:1901.10995*, 2019.
- [23] Lasse Espeholt, Hubert Soyer, Rémi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA: scalable distributed deep-rl with importance weighted actor-learner architectures. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1406–1415. PMLR, 2018.
- [24] Meng Fang, Tianyi Zhou, Yali Du, Lei Han, and Zhengyou Zhang. Curriculum-guided hindsight experience replay. In *NeurIPS*, 2019.
- [25] Yannis Flet-Berliac, Johan Ferret, Olivier Pietquin, Philippe Preux, and Matthieu Geist. Adversarially guided actor-critic. *CoRR*, abs/2102.04376, 2021.
- [26] Carlos Florensa, David Held, Markus Wulfmeier, Michael Zhang, and Pieter Abbeel. Reverse curriculum generation for reinforcement learning. In *Conference on robot learning*, pages 482–495. PMLR, 2017.
- [27] Sébastien Forestier, Rémy Portelas, Yoan Mollard, and Pierre-Yves Oudeyer. Intrinsically motivated goal exploration processes with automatic curriculum learning. *arXiv preprint arXiv:1708.02190*, 2017.
- [28] Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, et al. Threedworld: A platform for interactive multi-modal physical simulation. *arXiv preprint arXiv:2007.04954*, 2020.
- [29] Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.
- [30] Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [32] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [33] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In *International Conference on Machine Learning*, pages 4940–4950. PMLR, 2021.
- [34] Arthur Juliani, Ahmed Khalifa, Vincent-Pierre Berges, Jonathan Harper, Hunter Henry, Adam Crespi, Julian Togelius, and Danny Lange. Obstacle tower: A generalization challenge in vision, control, and planning. *CoRR*, abs/1902.01378, 2019.
- [35] Niels Justesen, Ruben Rodriguez Torrado, Philip Bontrager, Ahmed Khalifa, Julian Togelius, and Sebastian Risi. Procedural level generation improves generality of deep reinforcement learning. *CoRR*, abs/1806.10729, 2018.
- [36] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. In *Machine Learning*, pages 209–232. Morgan Kaufmann, 2002.
- [37] Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. ViZDoom: A Doom-based AI research platform for visual reinforcement learning. In *IEEE Conference on Computational Intelligence and Games*, pages 341–348, Santorini, Greece, Sep 2016. IEEE. The best paper award.
- [38] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of generalisation in deep reinforcement learning. *CoRR*, abs/2111.09794, 2021.

- [39] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of generalisation in deep reinforcement learning. *CoRR*, abs/2111.09794, 2021.
- [40] Heinrich Küttler, Nantas Nardelli, Thibaut Lavril, Marco Selvatici, Viswanath Sivakumar, Tim Rocktäschel, and Edward Grefenstette. Torchbeast: A pytorch platform for distributed RL. *CoRR*, abs/1910.03552, 2019.
- [41] Heinrich Küttler, Nantas Nardelli, Alexander H. Miller, Roberta Raileanu, Marco Selvatici, Edward Grefenstette, and Tim Rocktäschel. The nethack learning environment. *CoRR*, abs/2006.13760, 2020.
- [42] Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to personalized news article recommendation. In Michael Rappa, Paul Jones, Juliana Freire, and Soumen Chakrabarti, editors, *WWW*, pages 661–670. ACM, 2010.
- [43] Marlos C Machado, Marc G Bellemare, and Michael Bowling. Count-based exploration with the successor representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5125–5133, 2020.
- [44] Prasanta Chandra Mahalanobis. On the generalized distance in statistics. *Proceedings of the National Institute of Sciences (Calcutta)*, 2:49–55, 1936.
- [45] Jarryd Martin, Suraj Narayanan Sasikumar, Tom Everitt, and Marcus Hutter. Count-based exploration in feature space for reinforcement learning. *arXiv preprint arXiv:1706.08090*, 2017.
- [46] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.
- [47] Jesse Mu, Victor Zhong, Roberta Raileanu, Minqi Jiang, Noah D. Goodman, Tim Rocktäschel, and Edward Grefenstette. Improving intrinsic exploration with language abstractions. *CoRR*, abs/2202.08938, 2022.
- [48] Georg Ostrovski, Marc G Bellemare, Aäron Oord, and Rémi Munos. Count-based exploration with neural density models. In *International conference on machine learning*, pages 2721–2730. PMLR, 2017.
- [49] Pierre-Yves Oudeyer, Frdric Kaplan, and Verena V Hafner. Intrinsic motivation systems for autonomous mental development. *IEEE transactions on evolutionary computation*, 11(2):265–286, 2007.
- [50] Pierre-Yves Oudeyer and Frederic Kaplan. What is intrinsic motivation? a typology of computational approaches. *Frontiers in neurobotics*, 1:6, 2009.
- [51] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [52] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. *CoRR*, abs/1705.05363, 2017.
- [53] Aleksei Petrenko, Erik Wijmans, Brennan Shacklett, and Vladlen Koltun. Megaverse: Simulating embodied agents at one million experiences per second. In *International Conference on Machine Learning*, pages 8556–8566. PMLR, 2021.
- [54] Rémy Portelas, Cédric Colas, Katja Hofmann, and Pierre-Yves Oudeyer. Teacher algorithms for curriculum learning of deep rl in continuously parameterized environments. In *CoRL*, 2019.
- [55] Sébastien Racanière, Andrew Kyle Lampinen, Adam Santoro, David P. Reichert, Vlad Firoiu, and Timothy P. Lillicrap. Automated curricula through setter-solver interactions. *ArXiv*, abs/1909.12892, 2019.
- [56] Roberta Raileanu and Tim Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-generated environments. In *International Conference on Learning Representations*, 2020.

- [57] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [58] Mikayel Samvelyan, Robert Kirk, Vitaly Kurin, Jack Parker-Holder, Minqi Jiang, Eric Hambro, Fabio Petroni, Heinrich Küttler, Edward Grefenstette, and Tim Rocktäschel. Minihack the planet: A sandbox for open-ended reinforcement learning research. *CoRR*, abs/2109.13202, 2021.
- [59] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [60] Jürgen Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural controllers. In *Proc. of the international conference on simulation of adaptive behavior: From animals to animats*, pages 222–227, 1991.
- [61] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Claudia Pérez-D’Arpino, Shyamal Buch, Sanjana Srivastava, Lyne Tchapmi, et al. igibson 1.0: A simulation environment for interactive tasks in large realistic scenes. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7520–7527. IEEE, 2020.
- [62] J. Sherman and W. J. Morrison. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *Ann. Math. Stat.*, 21(1):124–127, 1950.
- [63] Bradley C Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- [64] Christopher Stanton and Jeff Clune. Curiosity search: Producing generalists by encouraging individuals to continually explore and acquire skills throughout their lifetime. *PLOS ONE*, 11(9):1–20, 09 2016.
- [65] Christopher Stanton and Jeff Clune. Deep curiosity search: Intra-life exploration improves performance on challenging deep reinforcement learning problems. *CoRR*, abs/1806.00553, 2018.
- [66] Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- [67] Sainbayar Sukhbaatar, Ilya Kostrikov, Arthur D. Szlam, and Rob Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. *ArXiv*, abs/1703.05407, 2018.
- [68] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [69] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [70] Haoran Tang, Rein Houthoofd, 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. *Advances in neural information processing systems*, 30, 2017.
- [71] Erik Wijmans, Abhishek Kadian, Ari S. Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *ICLR*, 2020.
- [72] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020.

- [73] Daochen Zha, Wenye Ma, Lei Yuan, Xia Hu, and Ji Liu. Rank the episodes: A simple approach for exploration in procedurally-generated environments. *arXiv preprint arXiv:2101.08152*, 2021.
- [74] Amy Zhang, Nicolas Ballas, and Joelle Pineau. A dissection of overfitting and generalization in continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*, 2018.
- [75] Chiyuan Zhang, Oriol Vinyals, Rémi Munos, and Samy Bengio. A study on overfitting in deep reinforcement learning. *CoRR*, abs/1804.06893, 2018.
- [76] Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E. Gonzalez, and Yuandong Tian. Noveld: A simple yet effective exploration criterion. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] We discuss how our method (as well as existing methods) are not able to solve MiniHack tasks with large action spaces in Section 5 as well as Appendix C.1.4.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Appendix A.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We included the code with instructions in the supplement, and will release the code on GitHub upon acceptance.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendix C.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We use the `rliable` library to compute error bars using stratified bootstrapping in Figure 5.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix C.1.5.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We cited all the codebases we built upon, see Section C.4.
 - (b) Did you mention the license of the assets? [Yes] We mentioned all licenses associated with the codebases we used, see Section C.4.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] All of our new code is included in our code release.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Broader Impact Statement

This work proposes an RL exploration method for contextual MDPs, which is a very broad framework. Many decision-making problems can be framed as contextual MDPs, such as autonomous driving (contexts represent cars/roads), household robotics (contexts represent houses), healthcare applications (contexts represent patients) and online recommendation/ad optimization (contexts represent customers). Like other RL exploration algorithms, our method facilitates learning a policy which maximizes some reward function specified by a designer. Depending on the goals of the reward function designer, executing the resulting policy could result in positive or negative consequences.

B Algorithm Details

Algorithm 2 Exploration via Episodic Elliptical Bonuses (E3B)

Initialize policy π , feature encoder ϕ and inverse dynamics model f .

while not converged **do**

 Sample context $c \sim \mu_C$ and initial state $s_0 \sim \mu_S(\cdot|c)$

 Initialize inverse covariance matrix: $C_0^{-1} = \frac{1}{\lambda}I$

for $t = 0, \dots, T$ **do**

$a_t \sim \pi(\cdot|s_t)$ // Sample action

$s_{t+1}, r_{t+1} \sim P(\cdot|s_t, a_t)$ // Step through environment

$b_{t+1} = \phi(s_{t+1})^\top C_t^{-1} \phi(s_{t+1})$ // Compute bonus

$u = C_t^{-1} \phi(s_{t+1})$

$C_{t+1}^{-1} = C_t^{-1} - \frac{1}{1+b_{t+1}} uu^\top$ // Update inverse covariance matrix

$\bar{r}_{t+1} = r_{t+1} + \beta b_{t+1}$

end for

 Perform policy gradient update on π using rewards $\bar{r}_1, \dots, \bar{r}_T$.

 Update ϕ and g using $\{(s_t, a_t, s_{t+1})\}_{t=0}^{T-1}$ to minimize the loss:

$$\ell = -\log(p(a_t | f(\phi(s_t), \phi(s_{t+1}))))$$

end while

The full algorithm details are shown above.

Additional intuition: The elliptical bonus is related to the Mahalanobis distance [44] which uses a similar bilinear form. However, the Mahalanobis distance would normalize the matrix C_{t-1} in equation 4.1 by the number of observations $t - 1$ in the episode, whereas the elliptical bonus does not. The elliptical bonus thus tends to decrease with the number of observations, similarly to the count-based bonus.

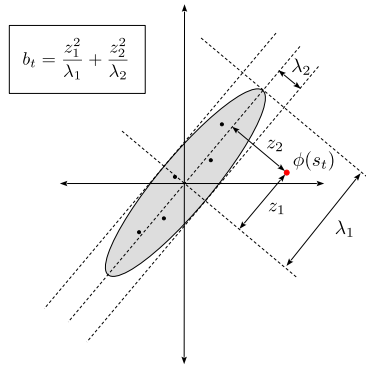


Figure 10: Illustration of the elliptical bonus in 2 dimensions.

C Experiment Details

C.1 MiniHack

C.1.1 Architecture Details

We follow the policy network architecture described in [58]. The policy network has four trunks: i) a 5-layer convolutional trunk which maps the full symbol image (of size 79×21) to a hidden representation, ii) a second 5-layer convolutional trunk which maps a 9×9 crop centered at the agent to a hidden representation, iii) an MLP trunk which maps the stats vector to a hidden representation, and iv) a 1-D convolutional trunk with interleaved max-pooling layers, followed by a fully-connected network which maps the message to a hidden representation. The hidden representations are then concatenated together, passed through a 2-layer fully-connected network followed by an LSTM [32] layer. The output of the LSTM layer is then passed to linear layers which produce action probabilities and a value function estimate.

The convolutional trunks i) and ii) have the following hyperparameters: 5 layers, filter size 3, symbol embedding dimension 64, stride 1, filter number 16 at each layer except the last, which is 8, and ELU non-linearities [16]. The MLP trunk iii) has 2 hidden layers of 64 hidden units each with ReLU non-linearities. The trunk iv) for processing messages has 6 convolutional layers, each with 64 input and output feature maps. The first two have kernel size 7 and the rest have kernel size 3. All have stride 1 and there are max-pooling layers (kernel size 3, stride 3) after the 1st, 2nd and 6th convolutional layers. The last two layers are fully-connected and have 128 hidden units and ReLU non-linearities.

For E3B, we used the same architecture as the policy encoder for the feature embedding ϕ , except we removed the last layers mapping the hidden representation to the actions and value estimate. The inverse dynamics model is a single-layer fully-connected network with 256 hidden units, mapping two concatenated ϕ outputs to a softmax distribution over actions.

C.1.2 RL Hyperparameters

For all algorithms we use IMPALA [23] as our base policy optimizer. Hyperparameters which are common to all methods are shown in Table 2. All algorithms were trained for 50 million environment steps. We did not anneal learning rates for any of the methods during training, since we found this yielded similar or better performance and simplified the setup.

Hyperparameters specific to E3B, NovelD, RIDE and ICM are shown in Tables 3, 4, 5 and 6. For both E3B and NovelD, we experimented with a rolling normalization of the intrinsic reward similar to that proposed in the RND paper [13]. Specifically, we maintained a running standard deviation σ of the intrinsic rewards and divided the intrinsic rewards by σ before feeding them to the policy optimizer. We found that this led to improved aggregate performance across environments for E3B and NovelD (the improvement for NovelD was relatively minor). We found that tuning the intrinsic reward coefficient was important for best performance for all methods, as well as the regularization coefficient of the C_t matrix λ for E3B. For RIDE, we used the default hyperparameters from the RIDE implementation in the MiniHack paper [58] and otherwise tuned the intrinsic reward coefficient. For ICM, we used the same hyperparameters for the forward and inverse models as for RIDE and also tuned the intrinsic reward coefficient.

Table 2: Common IMPALA Hyperparameters for MiniHack

Learning Rate	0.0001
RMSProp smoothing constant	0.99
RMSProp momentum	0
RMSProp ϵ	10^{-5}
Unroll Length	80
Number of buffers	80
Number of learner threads	4
Number of actor threads	256
Max gradient norm	40
Entropy Cost	0.0005
Baseline Cost	0.5
Discounting Factor	0.99

Table 3: Hyperparameters for E3B

Hyperparameter	Values considered	Final Value
Running intrinsic reward normalization	{True, False}	True
Ridge regularizer λ	{1.0, 0.1, 0.01}	0.1
Entropy Cost	{0.0005, 0.005}	0.005
Intrinsic reward coefficient β	{0.0001, 0.001, 0.01, 0.1, 1, 10}	1

Table 4: Hyperparameters for NovelD

Hyperparameter	Values considered	Final Value
Running intrinsic reward normalization	{True, False}	True
Scaling factor α	{0.1, 0.5}	0.1
Entropy Cost	{0.0005, 0.005}	0.005
Intrinsic reward coefficient β	{0.001, 0.01, 0.1, 1, 10, 100}	1

Table 5: Hyperparameters for RIDE

Hyperparameter	Values considered	Final Value
Forward Model loss coefficient	1.0	1.0
Inverse Model loss coefficient	0.1	0.1
Entropy Cost	{0.0005, 0.005}	0.0005
Intrinsic reward coefficient β	{0.001, 0.01, 0.1, 1, 10, 100}	0.1

Table 6: Hyperparameters for ICM

Hyperparameter	Values considered	Final Value
Forward Model loss coefficient	1.0	1.0
Inverse Model loss coefficient	0.1	0.1
Entropy Cost	{0.0005, 0.005}	0.0005
Intrinsic reward coefficient β	{0.001, 0.01, 0.1, 1, 10, 100}	0.1

C.1.3 Environment Details

We used 16 MiniHack environments in total. These include the following 9 navigation-based tasks:

```
'MiniHack-MultiRoom-N4-Locked-v0', 'MiniHack-MultiRoom-N6-Lava-v0',  
'MiniHack-MultiRoom-N6-Lava-OpenDoor-v0', 'MiniHack-MultiRoom-N6-LavaMonsters-v0',  
'MiniHack-MultiRoom-N10-OpenDoor-v0', 'MiniHack-MultiRoom-N10-Lava-OpenDoor-v0',  
'MiniHack-LavaCrossingS19N13-v0', 'MiniHack-LavaCrossingS19N17-v0',  
'MiniHack-Labyrinth-Big-v0'
```

as well as the following 7 skill-based tasks:

```
'MiniHack-Levitate-Potion-Restricted-v0', 'MiniHack-Levitate-Boots-Restricted-v0',  
'MiniHack-Freeze-Horn-Restricted-v0', 'MiniHack-Freeze-Wand-Restricted-v0',  
'MiniHack-Freeze-Random-Restricted-v0', 'MiniHack-LavaCross-Restricted-v0',  
'MiniHack-WoD-Hard-Restricted-v0'
```

The MultiRoom-N4-Locked and MultiRoom-N*-Lava, Labyrinth-Big, LavaCrossingS*N*, as well as all the skill-based tasks, are taken from the official MiniHack github repository (<https://github.com/facebookresearch/minihack>). We made some of these harder by increasing the number of rooms.

We also include some new environments which we designed to better test the limits of the different algorithms. Specifically, 'MiniHack-MultiRoom-N*OpenDoor-v0' are variants on the standard MultiRoom tasks, the only difference being that the doors connecting the rooms are initialized to be open rather than closed. This is because opening a door causes a message to appear "the door opens". By initializing the door to be closed, the agent does not see any messages when passing from one room to the next. As discussed in Appendix D, this seemingly trivial change can cause the NovelD-message variant to fail completely, shedding light on its lack of robustness.

We found that the MiniHack-MultiRoom-N6-Extreme task was impossible to solve consistently even for a human player due to the large quantities of monsters, so we did not include it. We instead included a variant with less monsters and lava called MultiRoom-N6-LavaMonsters.

C.1.4 Action Space Restriction

In their default setup, the skill-based tasks have a large context-dependent action space of 78 actions. Many of these actions are unrelated to the task at hand, but produce unique in-game messages when executed which nevertheless do not affect the underlying state of the MDP. For example, trying to pay a non-existent shopkeeper by executing the PAY action results in the message "There appears to be no shopkeeper here to receive your payment.", and executing the FIGHT action in the absence of adversaries results in the message "You attack thin air." We found that neither the baselines nor our proposed method were able to solve the skill-based tasks with the full action space. For the NovelD variants which use positions or symbolic images for the episodic counts, this is reasonable since the tasks are not navigation-based (these methods did not work even with restricted action spaces). For NovelD with the message-based bonus, we found that the agent ends up learning a policy where it executes many different actions which produce different messages, but do not change the underlying state of the game (such as the PAY and FIGHT actions described above). This makes sense from the perspective of optimizing intrinsic reward, since each new message seen within the episode provides additional intrinsic reward; however, this does not help explore the state space in a way that is helpful for discovering the true environment reward. We observed similar behavior for E3B, and hypothesize that the encoding learned through the inverse dynamics model keeps message information since this is useful for predicting actions, even if they have no real effect (such as PAY and FIGHT). In this case, the policy can then maximize intrinsic reward by executing these actions.

To avoid this unwanted behavior (which applies to all the methods we tested), we restricted the action space to only the actions which were useful for the tasks at hand. These included actions corresponding to direction movements, as well as different combinations of PICKUP, QUAFF, ZAP, FIRE, and WEAR (full details can be found in our code release). We denote the versions of the tasks with action space restriction with the "-Restricted-v0" suffix, as opposed to the original "-v0" suffix. We believe that better understanding this failure mode and designing exploration bonuses and/or representation learning methods which are robust to it is an important direction for future work.

C.1.5 Compute Details

Each algorithm was trained using 40 Intel(R) Xeon(R) CPU cores (E5-2698 v4 @ 2.20GHz) and one NVIDIA GP100 GPU. We used PyTorch [51] for all our experiments. Each run took between approximately 10 and 30 hours to complete. The total runtime depended on two factors: the computation time of the algorithm, and the behavior of the policy. In terms of algorithm computation time, E3B was roughly 1.5 to 2 times slower than the other methods, due to the fact that in our implementation an additional forward pass through the embedding network (used to compute the elliptical bonus) was performed on CPU. It is possible that a different implementation where this step would be done on GPU would be faster.

A second factor which influenced the total runtime was how quickly the agent learned to avoid dying. For certain environments (for example, those which contain lava), the agent can die quickly which causes the environment to be regenerated. For environments which call the MiniGrid library on the backend, this can be a performance bottleneck since generating new MiniGrid environments can be slow. We found that algorithms which cause the agent to die frequently were much slower on the `MiniHack-MultiRoom-N*Lava*` and `MiniHack-LavaCrossing*` environments, both of which are based on MiniGrid.

C.2 VizDoom

C.2.1 Architecture Details

For all VizDoom experiments, we used the same policy network architecture as in [56]: four 2D convolutional layers with 32 input and 32 output channels, kernel size 3×3 , stride 2 and padding 1, interleaved with Exponential Linear Unit (ELU) non-linearities [16]. The output of the convolutional layers feeds into a 2-layer LSTM [32] with 256 hidden units, followed by linear layers mapping the hidden units to a distribution over actions and a scalar value estimate.

The architecture for E3B’s encoder is identical to that of the policy network, except that the LSTM is replaced by a linear layer (there is thus no recurrence in the encoder network).

C.2.2 RL Hyperparameters

For our VizDoom experiments, we also used IMPALA as our base policy optimizer. For E3B we used the same IMPALA hyperparameters as for MiniHack. The hyperparameters specific to E3B, ICM and RIDE are given below. For both ICM and RIDE, we found that performance was quite sensitive to the coefficients of the forward and inverse dynamics model losses. We ended up using the hyperparameters from [56], which we found gave the best performance.

Table 7: Hyperparameters for E3B on VizDoom

Hyperparameter	Values considered	Final Value
Running intrinsic reward normalization	{False}	False
Ridge regularizer λ	{0.1, 0.01}	0.1
Intrinsic reward coefficient β	{ $3 \cdot 10^{-7}$, 10^{-6} , $3 \cdot 10^{-6}$, 10^{-5} , $3 \cdot 10^{-5}$, 10^{-4} }	$3 \cdot 10^{-6}$

Table 8: Hyperparameters for RIDE on VizDoom

Hyperparameter	Values considered	Final Value
Running intrinsic reward normalization	{False}	False
Forward loss coefficient	{0.5, 1.0}	0.5
Inverse loss coefficient	{0.1, 0.8}	0.8
Intrinsic reward coefficient β	{ 10^{-2} , $3 \cdot 10^{-3}$, 10^{-3} , $3 \cdot 10^{-4}$, 10^{-4} }	$3 \cdot 10^{-2}$

Table 9: Hyperparameters for ICM on VizDoom

Hyperparameter	Values considered	Final Value
Running intrinsic reward normalization	{False}	False
Forward loss coefficient	{0.2, 1.0}	0.2
Inverse loss coefficient	{0.1, 0.8}	0.8
Entropy cost	{0.0001, 0.005}	0.005
Intrinsic reward coefficient β	{ 10^{-2} , $3 \cdot 10^{-3}$, 10^{-3} , $3 \cdot 10^{-4}$, 10^{-4} }	$3 \cdot 10^{-3}$

C.3 Habitat

C.3.1 Environment Details

We used the HM3D [57] dataset, which consists of 1000 high-quality renderings of indoor scenes. Observations consist of 4 modalities: an RGB and depth image (shown in Figure 11a), GPS coordinates and the compass heading. The action space consists of 4 actions: $\mathcal{A} = \{\text{stop_episode}, \text{move_forward} (0.25\text{m}), \text{turn_left} (10^\circ), \text{turn_right} (10^\circ)\}$. The dataset scenes are split into 800/100/100 train/validation/test splits. Since the test split is not publicly available, we evaluate all models on the validation split. Each scene corresponds to a different context $c \in \mathcal{C}$ in the CMDP framework.

To measure exploration coverage, we compute the area revealed by the agent’s line of site using the function provided by the Habitat codebase ², which uses a modified version of Bresenham’s line cover algorithm. We define the exploration coverage to be:

$$\text{coverage} = \frac{\text{revealed area}}{\text{total area}}$$

See Figure 11b) for an illustration. For the results in Figure 8, we evaluated exploration performance for each algorithm by measuring its coverage on 100 episodes using scenes from the validation set (which were not used for training).



Figure 11: a) Visual observations in Habitat b) Exploration is measured as the proportion of the environment revealed by the agent’s line of sight over the course of the episode.

C.3.2 Architecture Details

For all Habitat experiments we used the same policy network as in [71], which includes a ResNet50 visual encoder [31] and a 2-layer LSTM [32] policy. In addition to RGB and Depth images, the agent also receives GPS coordinates and compass orientation, represented by 3 scalars total, which are fed into the policy. See the official code release at https://github.com/facebookresearch/habitat-lab/tree/main/habitat_baselines for full details.

For exploration algorithms which use inverse dynamics models (E3B and ICM), we set the architecture of the encoder ϕ to be identical to that of the policy network, except that the last layer mapping hidden units to actions is removed. The inverse dynamics model was a single layer MLP with 256 hidden units and ReLU non-linearities.

For exploration algorithms which use random network distillation (RND and NovelD), we set the architecture of the random network to be identical to that of the policy network.

C.3.3 RL Hyperparameters

The DD-PPO hyperparameters which are common to all the algorithms are listed in Table 10. The hyperparameters which are specific to each algorithm are listed in Table 11, 12, 4. For NovelD’s

²https://github.com/facebookresearch/habitat-lab/blob/main/habitat/visualizations/fog_of_war.py

count-based bonus, hashing the full image was too slow to be practical, so we subsampled images by a factor of 1000 used that for the count-based bonus, along with the GPS coordinates and compass direction.

Table 10: Common PPO/DD-PPO Hyperparameters for Habitat

Clipping	0.2
PPO epochs	2
Number of minibatches	2
Value loss coefficient	0.5
Entropy coefficient	0.00005
Learning rate	0.00025
ϵ	10^{-5}
Max gradient norm	0.2
Rollout steps	128
Use GAE	True
γ	0.99
τ	0.95
Use linear clip decay	False
Use linear LR decay	False
Use normalized advantage	False
Hidden size	512
DD-PPO Sync fraction	0.6

Table 11: Hyperparameters for E3B on Habitat

Hyperparameter	Values considered	Final Value
Ridge regularizer λ	{0.1}	0.1
Intrinsic reward coefficient β	{1.0, 0.1, 0.01, 0.001, 0.0001}	0.1
Inverse Dynamics Model updates per PPO epoch	3	3

Table 12: Hyperparameters for RND on Habitat

Hyperparameter	Values considered	Final Value
Intrinsic reward coefficient β	{1.0, 0.1, 0.01, 0.001, 0.0001}	0.1
Predictor Model updates per PPO epoch	3	3

Table 13: Hyperparameters for NovelD on Habitat

Hyperparameter	Values considered	Final Value
Intrinsic reward coefficient β	{1.0, 0.1, 0.01, 0.001, 0.0001}	0.1
Predictor Model updates per PPO epoch	3	3
Scaling factor α	0.1	0.1

Table 14: Hyperparameters for ICM on Habitat

Hyperparameter	Values considered	Final Value
Intrinsic reward coefficient β	{1.0, 0.1, 0.01, 0.001, 0.0001}	0.1
Forward Dynamics Model loss coefficient	{1.0}	1.0

C.3.4 Compute Details

Each job was run for 225 million steps, which took approximately 3 days on 32 GPUs with 10 CPU threads.

C.4 Codebases Used

Our codebase was built atop the following codebases:

- The official NovelD codebase: <https://github.com/tianjunz/NovelD> (Creative Commons Attribution-NonCommercial 4.0 license) for NovelD, RND, RIDE and count-based baselines (this codebase is build atop the official RIDE codebase below)
- The official MiniHack codebase: <https://github.com/facebookresearch/minihack> for network architectures appropriate to MiniHack environments (Apache 2.0 license)
- The official RIDE codebase: <https://github.com/facebookresearch/impact-driven-exploration> (Creative Commons Attribution-NonCommercial 4.0 license) for network architectures appropriate to VizDoom environments
- The official Habitat codebase: https://github.com/facebookresearch/habitat-lab/tree/main/habitat_baselines for Habitat experiments

D Additional Results and Discussion

D.1 Additional MiniGrid Results

Here we provide results for RIDE, AGAC and NovelD on additional MiniGrid environments used in prior work, with and without their count-based episodic terms. For all algorithms, we used the code released by the authors with the default hyperparameters³. Results in Figure 12 confirm the trend in Figure 1: in all cases, removing the count-based episodic term results in a failure to learn or makes learning much slower (such as for RIDE in MiniGrid-ObstructedMaze-2D1h-v0). We ran less seeds for AGAC because the official release’s multithreading implementation was not compatible with our cluster.

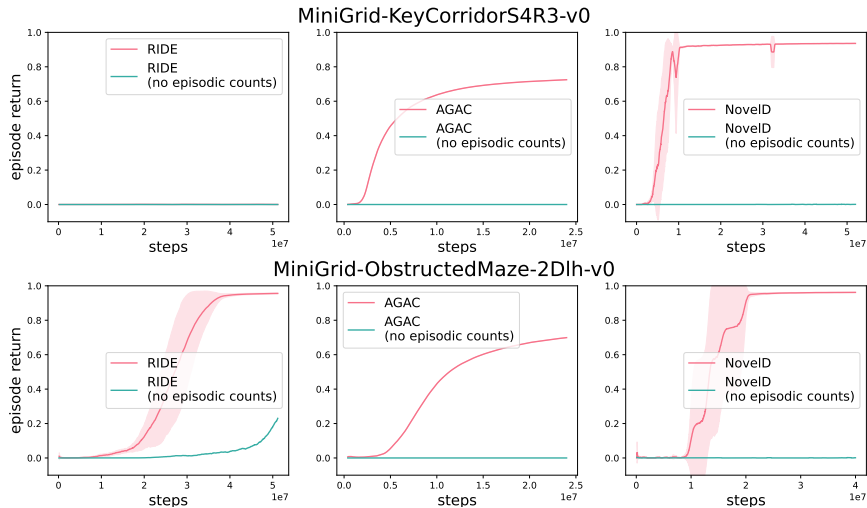


Figure 12: Results for RIDE, AGAC and NovelD with and without the count-based episodic bonus, over 5 random seeds (1 seed for AGAC). Shaded region indicates one standard deviation.

D.2 Effect of rank-1 updates

Here we compare the wall-clock time for computing the elliptical bonus using the rank-1 updates and the naive method of maintaining and inverting the full covariance matrix. We performed these experiments on the MiniHack-Freeze-Random-Restricted-v0 environment using the standard

³NovelD: <https://github.com/tianjunz/NovelD>, RIDE: <https://github.com/facebookresearch/impact-driven-exploration>, AGAC: <https://github.com/yfletberliac/adversarially-guided-actor-critic>

hyperparameter setup, on a machine with 80 Intel(R) Xeon(R) CPU cores (E5-2698 v4 @ 2.20GHz) and one NVIDIA GP100 GPU:

E3B using rank-1 update	767 FPS
E3B using matrix inversion	256 FPS

This shows that using the rank-1 update is essential for fast performance.

D.3 Negative results for RIDE with modified count-based bonuses

We ran some preliminary experiments investigating modifications to RIDE similar to those for NovelD. Specifically, the modified RIDE reward bonus is:

$$b(s_t) = \|\phi(s_{t+1}) - \phi(s_t)\|_2 \cdot \frac{1}{\sqrt{N_e(\psi(s_{t+1}))}}$$

where ϕ is the embedding learned with the inverse dynamics model and $\psi(s_{t+1})$ extracts some aspect of the state s_{t+1} . We investigated a version where $\psi(s_{t+1})$ extracts the (x_{t+1}, y_{t+1}) position information from the state (this method we called RIDE-POSITION and is the same version that was run in [58]) and a version where ψ extracts the message portion of the state (this method we called RIDE-MESSAGE). Despite tuning the intrinsic reward coefficient over the range $\{0.0001, 0.001, 0.01, 0.1, 1, 10\}$ on two tasks (MiniHack-MultiRoom-N10-v0 and MiniHack-Freeze-Horn-Restricted-v0) over 3 random seeds, none of the seeds was able to achieve positive reward for either of the tasks with either of the methods. Note that our results are consistent with those reported in [58], which also found that RIDE-POSITION did not outperform IMPALA on most tasks.

D.4 Additional MiniHack Results and Discussion

Results for all methods on individual tasks are shown in Figure 13. The first 9 tasks are navigation-based while the remaining 7 are skill-based.

First, as noted in the main text we see that when using the position bonus, NovelD succeeds in most of the tasks which primarily require the agent to spatially explore the environment, such as the Labyrinth-Big task or the different MultiRoom variants. On the other hand, this modality performs poorly on tasks based on skill acquisition, such as Freeze, LavaCross and WoD tasks. These tasks require the agent to pick up and use objects, and exploring the state space involves trying different action combinations which do not affect the agent’s position. Therefore, an exploration bonus which only encourages the agent to visit many different positions is not helpful.

The message bonus, on the other hand, succeeds on the skill acquisition tasks. On these tasks, the correct sequence of actions indeed causes a sequence of novel messages to appear. For example, let us consider the task Freeze-Wand. In this task, the agent is in a small room and must go to a wand, pick it up, pick the ZAP action, choose the wand and then choose a direction to zap. When conducting this sequence of actions, it encounters the following messages: "you see here a brass wand" (after navigating to the wand location), "f - a brass wand" (after executing the PICKUP action), "What do you want to zap? [f or ?*]" (after executing the ZAP action). Each time the agent visits a state-action pair required to solve the task, it receives one of these messages which provides it with a positive reward signal thanks to the message-based episodic bonus. This in turn reinforces the behavior leading to that state-action pair and allows it to ultimately solve the task.

When using the message bonus, NovelD fails completely on most of the navigation-based tasks, such as MultiRoom-N*-Lava-OpenDoor, LavaCrossingS19N13 and LavaCrossingS19N17. For these tasks, the agent must navigate its way through a series of rooms with lava walls or lava rivers with only one crossing. The only message it receives is "it’s a wall" if it runs into walls, which does not incentivize it to explore more locations which will eventually lead it to the goal. In contrast, the position bonus easily solves these tasks.

The symbolic image bonus performs poorly on the skill-based tasks, but performs well on some of the navigation-based ones. This is likely because for many of these tasks, the agent is the only

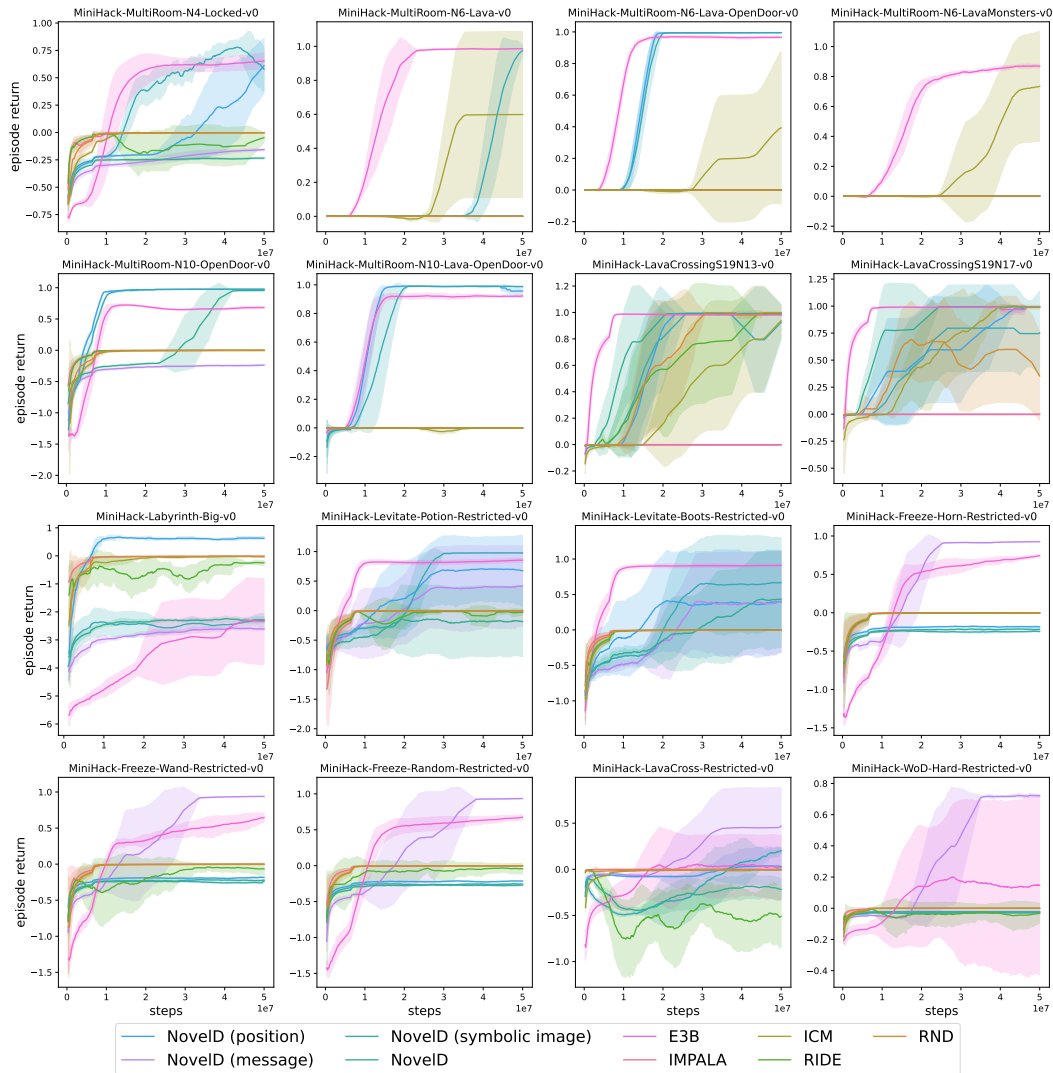


Figure 13: Mean results on individual tasks over 5 different seeds. Shaded region indicates one standard deviation.

moving entity, and hence the positional bonus and the image bonus will be similar. However, this is not always true. For example, in the Labyrinth-Big task, the environment is initially mostly hidden and gets revealed over time as the agent explores. This means that there is not a one-to-one correspondence between symbolic images and agent positions, as would be the case if the entire map was initially revealed. On this task, the position bonus succeeds but the symbolic bonus does not. Another particularly revealing example is the fact that the symbolic bonus succeeds on the MultiRoom-N6-Lava-OpenDoor environment, but fails on the MultiRoom-N6-Lava environment. The only difference between these two tasks is that in the former, all doors are initially open, which is marked by the - symbol. In the latter, they are initially closed (marked by the + symbol), and as they are opened by the agent, the symbol switches to -. In the open door version, there is a one-to-one correspondence between agent positions and symbolic images. On the closed door version, the number of possible symbols is the number of positions times the number of possible combinations of doors being open and closed—a much larger number. Once any door is opened, if the agent revisits a position it visited when the door was closed, it will once again receive bonus. This encourages the agent to revisit previously visited locations each time a door is opened, which does not align well with the task. This explains the poor performance of the symbolic image bonus on the closed door

version of the task. This again illustrates how count-based episodic bonuses can be sensitive to slight changes in the task’s construction.

When using the count-based bonus based on the full observation, which does not make any assumptions about which parts of it are useful for the task, NovelD fails on all the tasks except for two. A close look at the stats vector (see Figure 4) reveals that it contains a time counter which increments each time step: this effectively makes each observation unique, and hence the episodic bonus is constant⁴.

D.5 Sensitivity of λ regularizer

Figure 14 shows E3B’s performance on MiniHack for different values of the covariance regularizer λ . There is no statistically significant difference between $\lambda = 0.1$, which is the value we used in our experiments, and $\lambda = 0.01$ or $\lambda = 1.0$. For $\lambda = 0.001$, we observe a statistically significant drop in performance for the IQM metric only. This shows that E3B is fairly robust to the value of λ .

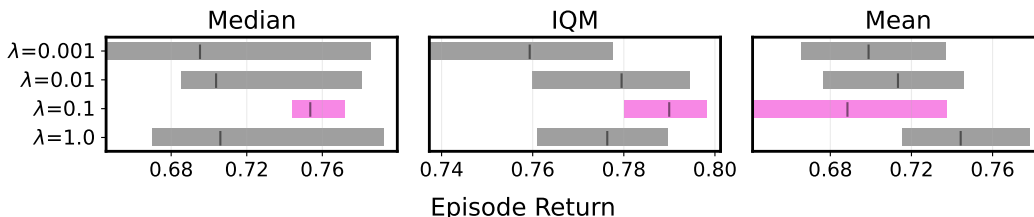


Figure 14: MiniHack performance for different values of λ parameter. All MiniHack experiments in the paper use the value $\lambda = 0.1$ unless otherwise noted. Intervals are computed using 5 random seeds using stratified bootstrapping.

⁴This is a simplification: there are some exceptions where the time counter is not incremented, such as if the agent runs into a wall or is searching through its inventory. This counts as a time step for the RL environment wrapper, but not for the underlying NetHack game engine. However these instances are relatively rare.