# Bayesian Curiosity for Efficient Exploration in Reinforcement Learning

Tom Blau[1], Lionel Ott[1], Fabio Ramos[1,2]

*Abstract*— Balancing exploration and exploitation is a fundamental part of reinforcement learning, yet most state-of-the-art algorithms use a naive exploration protocol like $\epsilon$-greedy. This contributes to the problem of high sample complexity, as the algorithm wastes effort by repeatedly visiting parts of the state space that have already been explored. We introduce a novel method based on Bayesian linear regression and latent space embedding to generate an intrinsic reward signal that encourages the learning agent to seek out unexplored parts of the state space. This method is computationally efficient, simple to implement, and can extend any state-of-the-art reinforcement learning algorithm. We evaluate the method on a range of algorithms and challenging control tasks, on both simulated and physical robots, demonstrating how the proposed method can significantly improve sample complexity.

## I. INTRODUCTION

Reinforcement Learning is a powerful approach for learning control policies for tasks that require a sequence of decisions, such as walking, playing a game, or navigating a maze. This learning is accomplished by repeatedly taking actions and observing both their effect on the environment state and some scalar reward signal that they elicit. The parameters of the policy are iteratively updated to increase the total expected reward over an entire sequence of observations and actions, called a trajectory. The weakness of RL is that it suffers from high sample complexity. Often millions of actions must be taken before an RL algorithm learns an efficient control policy, meaning a single execution can take days to complete. This problem is particularly pronounced in tasks where the space of possible actions and states is continuous and high-dimensional, a category that includes most real-world tasks, such as robotic object manipulation or autonomous driving.

One contributor to this high sample complexity is the exploration-exploitation trade-off [1]. This is a fundamental trade-off in any sequential decision making problem between *exploitation* - using available knowledge to follow a trajectory that maximizes expected reward - and *exploration* - searching previously unseen parts of the state space to find better trajectories. Excessive exploitation slows learning as effort is spent achieving high rewards without gaining new information that can be used to improve the policy. On the other hand, excessive exploration slows learning as effort is spent searching through regions of the state space where improved trajectories are unlikely to be found. Achieving good exploration is a challenging problem, as all but the simplest RL tasks have state spaces that are intractable to explore fully. The state space grows exponentially in the
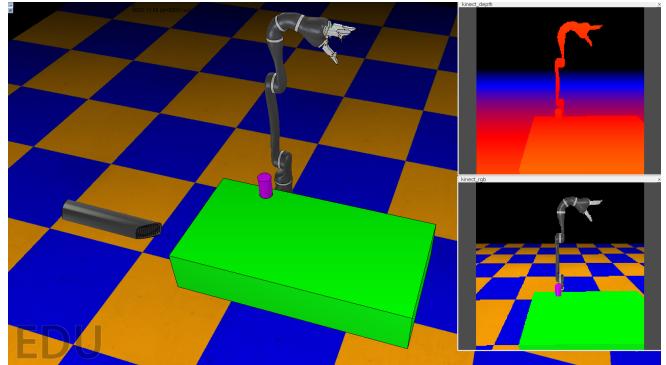


Fig. 1: Robotic reaching/grasping environment. On the left is a plain camera view of the environment, including the visual sensor. On the top right, the depth channel of the sensor displayed as a heatmap. On the bottom right, the RGB channels of the sensor.

number of dimensions, and the relationship between actions and state changes is often complex, so that even reaching a specific state is not trivial.

In spite of this, it is common for state-of-the-art RL algorithms to handle exploration in a naive manner. A standard approach is $\epsilon$-greedy, which consists of taking a random action with probability $\epsilon$ and following the learned policy with probability $1 - \epsilon$. The equivalent for continuous action spaces is adding Gaussian noise to policy actions. Consequently, exploration is often inefficient, spending much time searching through well-known regions of the state space that will provide little new information, or through regions that good trajectories are unlikely to intersect.

More recently, a class of algorithms has emerged that guides learning by use of intrinsic rewards [2], meaning a reward signal that comes from the model itself (intrinsic) rather than the environment (extrinsic). This is analogous to curiosity in humans, which are internally motivated to seek out novel experiences, even when there is no external reward for doing so. However, such algorithms tend to require cumbersome additions to the policy model.

In this paper we propose a curiosity based method that can be applied to an arbitrary RL algorithm in order to direct exploration towards parts of the state space that have not yet been visited. Using a Bayesian linear regression model with a learned latent space embedding, we compute the uncertainty of the model for arbitrary states. This uncertainty grows with the dissimilarity from previously seen points in the state space, and is therefore a good measure of the novelty of a given state observation. We use the model uncertainty to generate

* Correspondence to: Tom Blau, `tom.blau@sydney.edu.au`
[1] School of Computer Science, The University of Sydney, Australia
[2] NVIDIA, USA

an intrinsic reward signal that encourages visiting new states far from previously explored regions. Applying this approach to state-of-the-art RL algorithms, we verify experimentally that it accelerates learning in classic control tasks as well as in challenging robotics tasks with high dimensional state spaces, continuous action spaces, and even sparse rewards.

## II. RELATED WORK

Trading off exploration and exploitation is a fundamental problem in reinforcement learning, dating back to the multi-armed bandit problem [3] studied in statistics. The canonical algorithm for handling this trade-off was $\epsilon$-greedy [1], which consists of taking a random action with probability $\epsilon$ and otherwise taking a greedy action that maximizes the expected reward. Although guaranteed to converge, $\epsilon$-greedy is unable to incorporate information about the model's uncertainty. Many algorithms have been proposed to leverage such uncertainty information by estimating a posterior over the expected value, for example using Bayesian linear regression [4], [5] or Bayesian neural networks [6]. While provably efficient, these algorithms involve the use of an argmax operation, limiting them to problems with discrete actions. Our method seeks to address exploration in more realistic problems with continuous actions.

One powerful approach that can use uncertainty information to deal with the exploration-exploitation trade-off is the Upper Confidence Bound (UCB) algorithm [7]. Instead of greedily taking the action that maximizes the expected reward, the UCB algorithm takes the action that maximizes the sum of the reward's expectation and variance. A significant weakness, however, is that UCB does not extend naturally beyond the bandit setting. The EMU-Q algorithm [8] applies UCB principles to more general reinforcement learning problems. Bayesian linear regression with random Fourier features is used to estimate the Q-function as well as the uncertainty of the model and actions are taken that greedily maximize the sum of both. An alternative approach is SARSA with an uncertainty Bellman equation (UBE) [9]. The authors derive a different Bellman backup for the uncertainty of the Q-function, and approximate a Bayesian linear regression with a DQN [10]. Actions are taken by Thompson sampling [11]. While the EMU-Q algorithm works well for continuous action spaces with low-dimensional actions and observations, UBE works well in problems with high-dimensional observations but is restricted to discrete action spaces. In contrast, the method we propose can solve problems that have both a continuous action space and high-dimensional observations such as images.

In recent years, intrinsic motivation algorithms have emerged as a popular solution to the exploration-exploitation trade-off. This class of algorithms is based on the idea that exploration can be driven by a separate reward signal that encourages visiting under-explored states. Such a reward signal can be derived from visitation counts [12], [13], [14], model prediction error [15], [16], [17], variational information gain [18], or entropy maximization [19], [20]. In contrast, we propose to use the uncertainty of a Bayesian linear regression,

which is a well understood mathematical mechanism that can explicitly separate uncertainty in the model from uncertainty in the data. In [21], the uncertainty of a Gaussian process was used as a reward term in a linear-quadratic regularizer in order to encourage exploration. While Gaussian processes provide good uncertainty estimates, they are problematic when it comes to dealing with large data sets and high-dimensional data such as images, both of which our method can handle.

## III. BACKGROUND

### A. The Markov Decision Process

A Markov Decision Process (MDP) is described by a tuple $(\mathcal{S}, \mathcal{A}, \pi, \mathcal{P}, r, \gamma, H)$ where: $\mathcal{S}$ is the set of all possible states. $\mathcal{A}$ is the set of all possible actions. $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a stochastic policy mapping state-action pairs $(s, a)$ to the probability of choosing action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$. $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a state transition function mapping tuples $(s_t, a_t, s_{t+1})$ to the probability of arriving at state $s_{t+1}$ after taking action $a_t$ at state $s_t$. $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a reward function assigning a scalar reward value to each state-action pair. $\gamma \in (0, 1)$ is a discount factor while $H \in \mathbb{N}$ is a time horizon.

An MDP can be used to generate a sequences of states and actions as follows: given an initial state $s_0$, iteratively select the next action $a_t \sim \pi(s_t)$ and evolve the state by sampling $s_{t+1} \sim \mathcal{P}(s_t, a_t)$. The sequence $[s_0, a_0, \ldots, a_{H-2}, s_{H-1}]$ is called a *trajectory*. The expected total discounted reward over an entire trajectory is:

$$J_\pi = \mathbf{E}_\pi \left[ \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \right], \qquad (1)$$

where the expectation is subscripted by $\pi$ to denote the distribution over trajectories that $\pi$ induces, and $\gamma$ discounts rewards later in the trajectory. Let $\pi_\theta$ denote a parameterized policy whose parameters are $\theta$. A reinforcement learning algorithm seeks to maximize the expected total discounted reward by optimizing over the policy parameters. This can be expressed succinctly as the following optimization problem:

$$\theta^* = \arg\max_\theta J_{\pi_\theta}. \qquad (2)$$

### B. Bayesian Linear Regression

Linear regression is a form of regression analysis where the data is explained using a linear model [22]. That is, a model of the form $t = \mathbf{w} \cdot \mathbf{x} + b$, where $\mathbf{x}$ is the input, $\mathbf{w}$ is a vector of weights and $b$ is the intercept. Bayesian linear regression (BLR) [23] places a prior on the model weights $\mathbf{w}$. We choose a Gaussian prior:

$$\begin{aligned} p(\mathbf{w}) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0), \\ \mathbf{m}_0 &= \mathbf{0} \quad \text{and} \quad \mathbf{S}_0 = \alpha^{-1}\mathbf{I}, \end{aligned} \qquad (3)$$

where $\alpha$ is a hyperparameter representing the precision of the model, and $\mathbf{m_0}, \mathbf{S_0}$ are a vector and matrix describing the prior mean and covariance over the model weights. Note that in order to simplify the notation we remove the intercept $b$ and instead add a dummy dimension to the input $\mathbf{x}$ that will

always have a constant value of 1. Let there be a set of training data $(\mathbf{X}, \mathbf{t})$ where $\mathbf{X}$ is an $N \times D$ matrix whose rows are $D$-dimensional input vectors and $\mathbf{t}$ is an $N$-dimensional vector of targets. Although we restrict the targets to being scalar in order to simplify the notation, the following treatment can be easily extended to vector targets. In order to capture nonlinearity in the data, we transform the inputs of the BLR using a nonlinear function $\phi(x) : \mathbb{R}^D \mapsto \mathbb{R}^M$. Note that it is possible that $M \neq D$, meaning $\phi$ can change the dimensionality of the data. The posterior distribution over $\mathbf{w}$ given the training data is:

$$
\begin{aligned}
p(\mathbf{w}) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \\
\mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\boldsymbol{\Phi}^T\mathbf{t}), \\
\mathbf{S}_N &= \mathbf{S}_0^{-1} + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi},
\end{aligned} \quad (4)
$$

where $\beta$ is a hyperparameter representing noise precision in the data, $\mathbf{t}$ is a vector of targets and $\boldsymbol{\Phi}$ is an $N \times M$ matrix whose $i$-th row is $\phi(x_i)$. Note that after we perform an update with one set of data, we can perform another update with a second set by assigning $\mathbf{m}_0 = \mathbf{m}_N; \mathbf{S}_0 = \mathbf{S}_N$.

Linear regression gives a prediction of $t$ for an arbitrary input $\mathbf{x}$. BLR with the above formulation yields a posterior predictive distribution over $t$:

$$
\begin{aligned}
p(t|\mathbf{x}) &= \mathcal{N}(t|\mathbf{m}_N^T\phi(\mathbf{x}), \sigma_N^2(\mathbf{x})), \\
\sigma_N^2(\mathbf{x}) &= \beta^{-1} + \phi(\mathbf{x})^T\mathbf{S}_N\phi(\mathbf{x}),
\end{aligned} \quad (5)
$$

where $\sigma_N^2(\mathbf{x})$ is the variance of the predictive posterior.

## IV. BAYESIAN CURIOSITY

We propose *Bayesian Curiosity*, a method that can extend any RL algorithm by adding a secondary model that produces an intrinsic reward signal. The pipeline of the method can be seen in fig. 2. In black is the standard RL machinery: a state observation $o_t$ is passed to the policy network $\pi_\theta$, which produces action $a_t$. This action is then passed to the environment $E$, which evolves its state and emits both an environment reward $e_t$ and the next state observation. In green is the proposed Bayesian Curiosity mechanism: $o_t$ is given as input to a neural network $\phi_\psi$ with parameters $\psi$, which transforms it into a latent space. This latent space representation is passed to a BLR $B$, which computes the uncertainty for the given input. The uncertainty is then used to construct a curiosity reward signal $c_t$, which together with $e_t$ forms the reinforcement reward $r_t$.
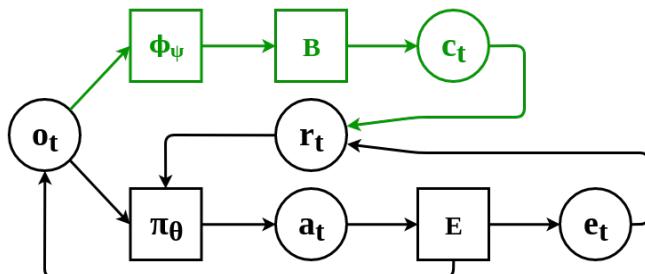


Fig. 2: Pipeline of RL with Bayesian Curiosity. Squares are functions and circles are variables.

### A. Latent Space Embedding

BLR is a linear model, and requires some transformation $\phi$ on the inputs in order to capture nonlinearity in the data. Further, it scales poorly with the dimensionality of the data. We will therefore learn a transformation using a neural network $\phi_\psi$ with weights $\psi$ to capture nonlinearity and reduce dimensionality. Algorithm 1 introduces a procedure for jointly optimizing the parameters of the network $\phi_\psi$ and BLR $B$, given a set of training data $D$. We acquire the training data by generating expert demonstrations with small Gaussian noise, meaning that our method is a form of learning from demonstration. In terms of the BLR, this means that $\mathbf{x}$ will be observations and $\mathbf{t}$ will be actions.

To train $\phi_\psi$ we follow a stochastic gradient descent procedure defined in algorithm 1. In each epoch, randomly sample a subset of the training data $E \subseteq D$. We use $E$ as the data for the BLR update in (4), denoting the parameters of the posterior over BLR weights as $\mathbf{m}_N^*$ and $\mathbf{S}_N^*$. Note that $\mathbf{m}_N^*$ and $\mathbf{S}_N^*$ are expressions parameterized by $\psi$. Plugging this into (5), we get the moments of the predictive posterior for an arbitrary demonstration $(\mathbf{x_i}, \mathbf{t_i}) \in D$:

$$
\begin{aligned}
\mu_N^*(\mathbf{x_i}) &= \mathbf{m}_N^{*T}\phi_\psi(\mathbf{x_i}), \\
\sigma_N^{*2}(\mathbf{x_i}) &= \beta^{-1} + \phi_\psi(\mathbf{x_i})^T\mathbf{S}_N^*\phi_\psi(\mathbf{x_i}).
\end{aligned} \quad (6)
$$

The loss for our SGD is the negative log-likelihood (NLL) of $B$ for $(\mathbf{x_i}, \mathbf{t_i})$, given by:

$$
L(x_i, t_i) = \frac{\log(2\pi)}{2} + \frac{\log(\sigma_N^{*2}(\mathbf{x_i}))}{2} + \frac{(t_i - \mu_N^*(\mathbf{x_i}))^2}{\sigma_N^{*2}(\mathbf{x_i})}. \quad (7)
$$

In each epoch, having chosen $E$, we repeatedly sample minibatches $mb \subset D$ without replacement, and take stochastic gradient steps to minimize the NLL loss w.r.t. $\psi$.

### B. Reinforcement Learning with Bayesian Curiosity

Given the latent space embedding $\phi_\psi$, we define the curiosity reward as:

$$
c_t = \log(\sigma^2(o_t)) = \log\left(\beta^{-1} + \phi_\psi(o_t)^T\mathbf{S}_N\phi_\psi(o_t)\right), \quad (8)
$$

where $\sigma^2$ is the variance of the predictive posterior in (5) and $o_t$ is the observation at time $t$. Since $\sigma^2$ is bounded from

---

**Algorithm 1** Learning Latent Embedding

**Input:** Training data $D$
Randomly initialize $\phi_\psi$
Initialize $\mathbf{m}_0, \mathbf{S}_0$ according to (3)
**repeat**
    Randomly select subset $E \subseteq D$
    Randomly partition $D$ into minibatches $MB$
    **for** $mb \in MB$ **do**
        Compute NLL for $mb$ and $E$ according to (7)
        Take SGD step w.r.t. $\psi$ to minimize NLL
    **end for**
**until** improvement on the NLL loss has converged

below by $\beta^{-1}$, the lower bound of $c_t$ is $-\log(\beta)$. Given extrinsic reward $e_t$, the final combined reward is:

$$r_t = e_t + \eta \cdot c_t, \qquad (9)$$

where $\eta$ is a hyperparameter that controls the weight of the curiosity reward.

We proceed to perform reinforcement learning with the combined reward signal $r_t$ as shown in algorithm 2. Let $A$ be some RL algorithm that can be chosen arbitrarily. Before RL begins, we reset the Bayesian linear regression $B$. In every episode, we execute rollouts of the policy $\pi_\theta$ to obtain sequences of actions, observations, and rewards. The observations $o_t$ are used to update $B$. Since $B$ is a BLR, uncertainty will, in expectation, be higher for new observations whose latent representations have low cosine similarity to previous observations. According to [9], states differing only in ways irrelevant to the task will be mapped to similar representations, and vice versa. Thus the curiosity reward will be higher for novel states. Note that we do not need to update $\mathbf{m}_N$, since we only use the uncertainty of $B$. This frees us from having to find target values $\mathbf{t}$, for which there is in general no intuitive candidate in the reinforcement learning setting. At the end of every episode $\pi_\theta$ is updated according to the rules of the algorithm $A$. Note that throughout the RL stage $\phi_\psi$ is not updated, as doing so would change the latent space and thus invalidate all previous updates to $B$.

## V. EXPERIMENTAL RESULTS

We proceed to test the effectiveness of our method on a series of continuous control and robotics tasks. In the following, all neural networks are implemented using Theano [24]. Pre-training is done using ADAM [25] with $l_2$ regularization. For the RL algorithms we use the implementations in RLLab [26]. Hyperparameters are $\alpha = 1e-4$ and $\beta = 1e2$ for the BLR, and the curiosity weight is $\eta = 1$. Source code is available at https://gitlab.com/tomblau/Bayesian-Curiosity.

---

**Algorithm 2** Bayesian Curiosity

---

**Input:** Environment $\mathcal{E}$, time horizon $H$, RL algorithm $A$, learned transformation $\phi_\psi$, untrained BLR $B$
Randomly initialize $\pi_\theta$
Initialize $\mathbf{S}_0$ according to (3)
**repeat**
    Reset $\mathcal{E}$ and observe $o_0$
    Create empty buffer $Z$
    **for** $i = 1$ **to** $H$ **do**
        $a_{i-1} = \pi_\theta(o_{i-1})$
        Execute $a_{i-1}$ and observe $o_i, e_i$
        Compute $r_i$ according to (8) and (9)
        Add $(a_{i-1}, o_{i-1}, r_i)$ to $Z$
    **end for**
    Update $\mathbf{S}_0$ according to (4) using $\phi_\psi(o_{0:H-1})$
    Update $\pi_\theta$ using buffer $Z$ and algorithm $A$
**until** policy improvement has converged

---

### A. Classic Control

We begin by evaluating our method on classic control problems with a number of RL algorithms: REINFORCE [27], DDPG [28], and TRPO [29] in four continuous control tasks: Mountaincar, Cartpole Swingup, Pendulum, and Acrobot. For each algorithm and environment we evaluate the algorithm with and without Bayesian Curiosity. As an additional exploration baseline, we also evaluate the baseline algorithm combined with VIME [18]. Extrinsic rewards for the above tasks have been sparsified to emphasize the importance of exploration.

Table I shows the results for this set of experiments, with each entry aggregated over 10 random seeds. The "reward" columns show the median total trajectory reward, as well as the lower and upper quartiles (the $25^{th}$ and $75^{th}$ percentiles, respectively). The "speedup" columns show how much more quickly the algorithms with Bayesian Curiosity learned compared with their respective baselines. The value is the ratio of the number of timesteps the baseline took to achieve its best result and the number of timesteps the corresponding Bayesian Curiosity algorithm took to match the baseline. Results for DDPG-VIME are missing as the authors of the original paper have not made code available for this algorithm. For most tasks and algorithms, our method is able to achieve comparable or superior performance in significantly fewer timesteps compared with both the standard RL and intrinsic motivation baselines. This is in spite of the fact that these are relatively simple problems with low-dimensional state and action spaces. Notable exceptions are DDPG on the Mountaincar task and TRPO on the Pendulum task, in which Bayesian Curiosity achieves lower performance than the baselines. In these cases, "speedup" is the ratio of the number of timesteps Bayesian Curiosity took to achieve its best result and the number of timesteps the baseline took to match it. For DDPG on Mountaincar, learning for both the baseline and Bayesian Curiosity versions is unstable, and the agents will periodically forget and re-learn how to achieve successful trajectories. Indeed, the lower quartile for both versions is $-200$ throughout the learning process, meaning that in every episode at least a quarter of agents have failed
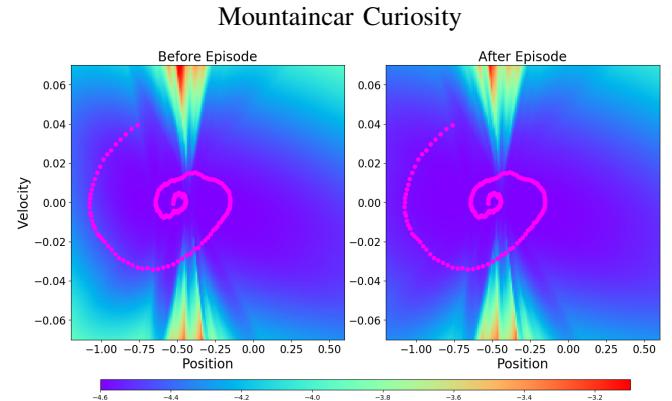


Fig. 3: Curiosity reward surface for mountaincar before and after an episode. Dots indicate visited states

TABLE I: Relative improvement of using Bayesian Curiosity in classic control tasks. Showing median and interquartile range of rewards over 10 experiments, as well as relative speedup in achieving peak performance.

| | | TRPO | | DDPG | | REINFORCE | |
|---|---|---|---|---|---|---|---|
| | | Reward | Speedup | Reward | Speedup | Reward | Speedup |
| Mountaincar | Curiosity | 20.0 (16.0, 22.0) | - | -50.0 (-200.0, -7.0) | - | **30.15** (4.80, 31.70) | - |
| | VIME | **22.0** (17.0, 24.0) | 0.8 | - | - | 28.5 (-15.75, 30.375) | 3.36 |
| | Vanilla | 18.0 (12.25, 22.0) | 1.43 | **11.0** (-200, 20.0) | 0.67 | 29.85 (4.75, 30.9) | 2.25 |
| Swingup | Curiosity | **-1.39** (-15.80, 12.36) | - | **33.11** (-17.78, 44.33) | - | **-46.53** (-57.61, -35.05) | - |
| | VIME | -46.66 (-67.93, -33.34) | 6.76 | - | - | -68.57 (-81.67, -56.04) | 5.33 |
| | Vanilla | -23.60 (-45.99, -9.94) | 4.44 | 25.0 (-20.71, 42.22) | 1.61 | -48.09 (-57.99, -37.47) | 1.11 |
| Pendulum | Curiosity | -80.33 (-87.39, -74.47) | - | **-27.13** (-50.2, -14.43) | - | **-80.53** (-86.22, -76.41) | - |
| | VIME | **-79.16** (-86.03, -72.03) | 0.833 | - | - | -81.43 (-87.11, -77.19) | 2.52 |
| | Vanilla | -79.63 (-88.25, -66.75) | 0.89 | -53.51 (-100, -19.45) | 13.33 | -83.91 (-100.0, -77.96) | 9.09 |
| Acrobot | Curiosity | **-182.7** (-241.33, -148.83) | - | **-138.05** (-203.25, -110.13) | - | **-135.5** (-190.08, -108.4) | - |
| | VIME | -230.1 (-319.83, -183.58) | 2.47 | - | - | -137.15 (-191.33, -110.40) | 1.14 |
| | Vanilla | -197.8 (-286.88, 143.05) | 1.45 | -143.1 (-208.8, -111.95) | 1.52 | -145.9 (-224.13, -115.0) | 2.21 |

to find a successful trajectory. For TRPO on Pendulum, the baseline and Bayesian Curiosity agents achieve similar results.

Finally, we examine how the surface of curiosity rewards changes when novel states are discovered. Figure 3 shows heatmaps of curiosity rewards w.r.t. the state space before and after an episode of mountaincar. The states visited in this episode are shown as pink dots in both plots. There is a noticeable decrease in curiosity in the $[-1.2, -0.75]$ range of the x-axis, where new states were visited that have not been seen before.

### B. Robotic Control

We now examine the effect of Bayesian Curiosity rewards on learning sensorimotor control policies for a 6-DOF robotic arm. This includes *Cartesian* control problems, where goal state information is represented by its 3-d Cartesian coordinates, and *visuomotor* control problems, which instead have image data. All tasks are executed in the V-REP simulation environment [30] using a Kinova Jaco arm, as shown in fig. 1, except for the final task which is executed on a real robot. Training data for the curiosity model are generated using an inverse kinematics solver, and the set consists of 8000 data points. To speed up the experiments and to make comparisons with the baselines more fair, all policies are pre-trained to imitate this training dataset. Experiments focus on TRPO because it was found to produce the most stable learning of all tested algorithms. All policies have Gaussian action noise determined by a neural network whose parameters are learned by the RL algorithm.

We begin with two Cartesian problems, reaching and grasping. The arm must be controlled to reach or grasp a cylinder that sits on a table. The cylinder is placed uniformly at random in a $40cm^2$ square at the beginning of an episode. In the reaching task, this cylinder is a non-interactable object that serves as an objective marker. Collision of the arm

with itself or other objects results in failure. Rewards are dense for the reaching task and sparse for the grasping task. Observations consist of the angles of the 6 controllable joints as well as the 3-dimensional coordinates of the cylinder. Actions are a vector of angle deltas, prescribing a change in the joint angles.

The top row of fig. 4 shows performance on the Cartesian reaching and grasping tasks for TRPO with and without Bayesian Curiosity. We also include VIME [18] and RND [15] as intrinsic motivation baselines. Performance is measured as $n_{succ}/n_{tsteps}$, where $n_{succ}$ is the number of successes and $n_{tsteps}$ is the number of timesteps in the 10 most recent episodes. This metric directly captures the ability of the policy to successfully complete the task, and unlike the success rate it also assigns a higher value to completing the task in fewer timesteps. In the reaching task, although median performance is similar, the lower quartile improves much more quickly in our method compared with the baselines. For the baselines, some initializations result in agents that take a very long time to converge, whereas with Bayesian Curiosity even the slowest agents converge quickly. This indicates that our method is more robust to random initialization.

For both the Cartesian and visuomotor cases, the grasping task exhibits a much starker difference between our method and the baselines than the reaching task does. The lower quartile of Bayesian Curiosity quickly rises above the *upper* quartile of the baselines. Further, whereas the trendlines for the baseline agents are almost linear, the curiosity agents exhibits a more sigmoidal trendline, improving rapidly at first before slowing down. There are two factors contributing to this difference. The first is that the grasping tasks have sparse rewards while the reaching tasks have dense rewards. The second has to do with the nature of the grasping task – any success state is very close to a failed state wherein the cylinder was knocked over rather than grasped. This
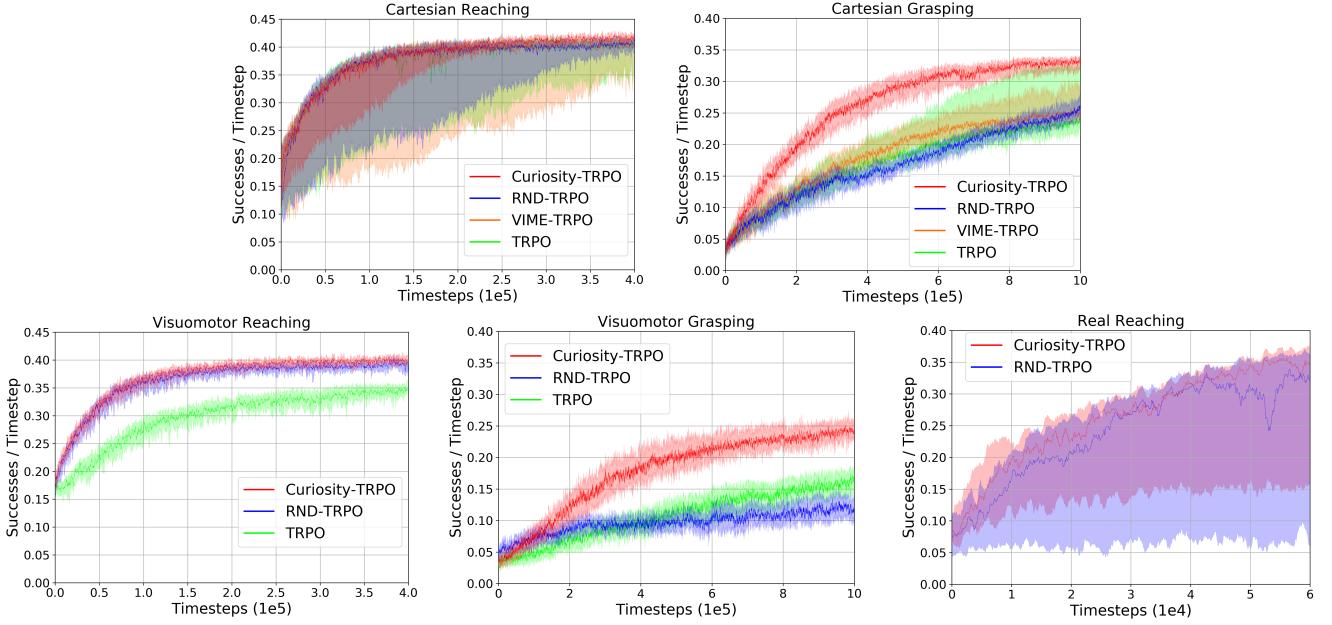
Fig. 4: Median and interquartile range for robotic reaching and grasping, comparing TRPO with Bayesian Curiosity, with RND, with VIME, and with only extrinsic rewards. Performance is measured as the ratio of successes to timesteps. Aggregated over 5 random seeds.

means that there is a bottleneck in the state space, formed by failure states, which must be traversed to reach a success state. This kind of geometry is very difficult to explore using Gaussian action noise, and comparatively easier to explore with a curiosity reward that discourages revisiting explored states.

In the next set of experiments we investigate visuomotor versions of the reaching and grasping tasks. Observations no longer include the 3-d coordinates of the cylinder, but instead contain RGBD images taken from a fixed camera pose. The performance of our method is compared with the vanilla TRPO and RND baselines (VIME does not scale to high-dimensional observations) in the bottom row of fig. 4. In the reaching task, our method greatly outperforms the naive baseline, quickly achieving performance that vanilla TRPO can't match even after 5 times as many timesteps, and slightly improves on the RND baseline. Compared with the Cartesian case, the intrinsic motivation methods perform better, while the performance of the naive baseline degrades. This suggests that agents with intrinsic motivation can explore enough to leverage the additional information of the visual sensors, while agents that rely on Gaussian noise for exploration struggle to explore the enlarged observation space. In the grasping task, our method achieves performance that the baselines can't match even after training for 3 times as many timesteps.

Finally, we replicate the Cartesian reaching task on a physical robotic arm. The bottom-right plot in fig. 4 shows the results of this set of experiments. Due to the time requirements of robotic experiments only RND, the strongest baseline in simulation, was included. As in the simulated reaching task, we see that median performance is comparable to the RND baseline, while the lower quartile grows more quickly,

indicating higher robustness to random initialization. This robustness is particularly valuable when dealing with physical robots, as each run of the algorithm is expensive and time-consuming.

## VI. CONCLUSIONS

In this work we introduced a new method that combines BLR with a learned latent space embedding to generate a curiosity signal that directs exploration towards novel states, and demonstrated its capability to augment a variety of standard RL algorithms. Compared with both naive Gaussian noise exploration and SOTA intrinsic motivation methods, our method is able to accelerate exploration and achieve comparable or superior performance in fewer timesteps. This improvement is particularly noticeable when the state space has a geometry that makes it difficult to explore, or when the reward function is sparse, and provides little information to help direct exploration. The ability to learn policies from sparse rewards is highly desirable, as it obviates the need for carefully designing reward functions using expert knowledge.

Future work can adapt Bayesian Curiosity to continually update the latent space embedding during the RL procedure, eliminating the need for demonstrations. Another possible extension is to combine the latent space representation with Random Fourier Features, which can approximate a Gaussian Process [31].

## REFERENCES

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press Cambridge, 1998.

[2] P. Y. Oudeyer and F. Kaplan, "How can we define intrinsic motivation?" in *International Conference on Epigenetic Robotics*, 2008.

[3] H. Robbins, "Some aspects of the sequential design of experiments," in *Herbert Robbins Selected Papers*. Springer, 1985.

[4] I. Osband, B. Van Roy, and Z. Wen, "Generalization and exploration via randomized value functions," in *International Conference on Machine Learning*, 2016.

[5] K. Azizzadenesheli, E. Brunskill, and A. Anandkumar, "Efficient exploration through bayesian deep q-networks," in *IEEE Information Theory and Applications Workshop*, 2018.

[6] Z. Lipton, X. Li, J. Gao, L. Li, F. Ahmed, and L. Deng, "Bbq-networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue systems," in *AAAI Conference on Artificial Intelligence*, 2018.

[7] P. Auer, "Using confidence bounds for exploitation-exploration trade-offs," *Journal of Machine Learning Research*, 2002.

[8] P. Morere and F. Ramos, "Bayesian rl for goal-only rewards," in *Conference on Robot Learning*, 2018.

[9] B. O'Donoghue, I. Osband, R. Munos, and V. Mnih, "The uncertainty bellman equation and exploration," in *International Conference on Machine Learning*, 2018.

[10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, B. M. G., A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, "Human-level control through deep reinforcement learning," *Nature*, 2015.

[11] W. R. Thompson, "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples," *Biometrika*, 1933.

[12] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos, "Unifying count-based exploration and intrinsic motivation," in *Advances in Neural Information Processing Systems*, 2016.

[13] M. C. Machado, M. G. Bellemare, and M. Bowling, "Count-based exploration with the successor representation," *arXiv preprint*, 2018.

[14] H. Tang, R. Houthooft, D. Foote, A. Stooke, *et al.*, "# exploration: A study of count-based exploration for deep reinforcement learning," in *NIPS*, 2017.

[15] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, "Exploration by random network distillation," *International Conference on Learning Representations*, 2019.

[16] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *International Conference on Machine Learning*, 2017.

[17] J. Achiam and S. Sastry, "Surprise-based intrinsic motivation for deep reinforcement learning," *arXiv preprint*, 2017.

[18] R. Houthooft, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel, "Vime: Variational information maximizing exploration," in *Advances in Neural Information Processing Systems*, 2016.

[19] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," in *ICLR*, 2019.

[20] Z. Hong, T. Shann, S. Su, Y. Chang, *et al.*, "Diversity-driven exploration strategy for deep reinforcement learning," in *NIPS*, 2018.

[21] S. Bechtle, A. Rai, Y. Lin, L. Righetti, and F. Meier, "Curious ilqr: Resolving uncertainty in model-based rl," in *ICML Workshop on Reinforcement Learning for Real Life*, 2019.

[22] X. Yan and X. Su, *Linear Regression Analysis: Theory and Computing*. World Scientific, 2009.

[23] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, 2006.

[24] J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley, and Y. Bengio, "Theano: a CPU and GPU math expression compiler," in *Python for Scientific Computing Conference (SciPy)*, 2010.

[25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint*, 2014.

[26] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," *International Conference on Machine Learning*, 2016.

[27] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Advances in Neural Information Processing Systems*, 2000.

[28] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint*, 2015.

[29] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust region policy optimization," *International Conference on Machine Learning*, 2015.

[30] E. Rohmer, S. P. N. Singh, and M. Freese, "V-rep: a versatile and scalable robot simulation framework," in *International Conference on Intelligent Robots and Systems*, 2013.

[31] J. Quinonero-Candela, C. E. Rasmussen, A. R. Figueiras-Vidal, *et al.*, "Sparse spectrum gaussian process regression," *Journal of Machine Learning Research*, vol. 11, 2010.