
Ready Policy One: World Building Through Active Learning

Philip Ball*
University of Oxford
ball@robots.ox.ac.uk

Jack Parker-Holder*
University of Oxford
jackph@robots.ox.ac.uk

Aldo Pacchiano
UC Berkeley
pacchiano@berkeley.edu

Krzysztof Choromanski
Google Brain Robotics
kchoro@google.com

Stephen Roberts
University of Oxford
sjrob@robots.ox.ac.uk

Abstract

Model-Based Reinforcement Learning (MBRL) offers a promising direction for sample efficient learning, often achieving state of the art results for continuous control tasks. However many existing MBRL methods rely on combining greedy policies with exploration heuristics, and even those which utilize principled exploration bonuses construct dual objectives in an ad hoc fashion. In this paper we introduce Ready Policy One (RP1), a framework that views MBRL as an active learning problem, where we aim to improve the world model in the fewest samples possible. RP1 achieves this by utilizing a hybrid objective function, which crucially adapts during optimization, allowing the algorithm to trade off reward v.s. exploration at different stages of learning. In addition, we introduce a principled mechanism to terminate sample collection once we have a rich enough trajectory batch to improve the model. We rigorously evaluate our method on a variety of continuous control tasks, and demonstrate statistically significant gains over existing approaches.

1 Introduction

Reinforcement Learning (RL) considers the problem of an agent learning to construct of actions that result in an agent receiving high rewards in a given environment. This can be achieved in various ways such as: learning an explicit mapping (*a policy*) from states to actions that maximizes expected return (policy gradients), or inferring such a mapping by calculating the expected return for a given state-action pair (TD-control methods). Model-Based Reinforcement Learning (MBRL) seeks to improve the above by learning a model of the dynamics (from agent’s interactions with the environment) that can be leveraged across many different tasks (transferability) and for planning, which is substantially less expensive than a real environment (which plays crucial role in robotics applications).

A series of recent work [18, 22, 31, 10, 9, 30, 44] illustrate the benefits of MBRL-approaches that allow us to decouple learning task-dependent policy and task-agnostic dynamics. With recent advances, MBRL approaches often outperform model-free methods [51]. However, these results are often overly sensitive to heuristics.

In particular, many of these methods lack a principled mechanism to acquire data for training the model. This issue is circumvented in [18], since they only consider environments which can be

*Equal contribution.

explored with random policies. Other model-based approaches, such as [31, 16, 21], rely on stochastic policies to aid exploration, and inevitably acquire redundant data which reduces sample efficiency. Such issues have been highlighted previously [42, 48], and motivate the design of our algorithm. Concretely, we reduce the cost incurred from data collection by using active learning methods, and introduce an early stopping mechanism to address the issue of redundancy.

Efficient exploration is a challenge for existing RL algorithms, and is a core focus in model-free RL [8, 25, 3, 32]. Despite often considering a principled objective, these methods generally contain a fixed temperature parameter, thus requiring hand engineering to determine the optimal degree of exploration. Our approach adjusts this parameter in an online manner from the collected trajectories, and we provide an information theoretic motivation for our exploration objective.

In this paper, we introduce a novel approach to acquiring data for training world models through exploration. Our algorithm, Ready Policy One (RP1), includes principled mechanisms which acquire data for model-based RL through the lens of Online Active Learning. Crucially, we continue to jointly optimize our policies for both reward and model uncertainty reduction, since we wish to avoid focusing on stochastic or challenging regions of the state space which have no impact on the task at hand. Therefore policies used for data collection also perform well in the true environment, and means we can use these policies for evaluation. Consequently, a separate ‘exploit’ agent does not need to be trained [23].

To summarize, our key contributions are as follows:

- Inspired by Active Learning, we train policies in a learned world model with the objective of acquiring data that most likely leads to subsequent improvement in the model.
- We introduce a novel early-stopping criteria for real-environment samples, reducing redundancy in expensive data collection.
- We adapt the objective function as learning progresses using an Online Learning mechanism.

The paper is structured as follows. In Section 2 we discuss related work. In Section 3 we describe our RL setting and introduce basic concepts. In Section 4 we introduce our method and related theory. Finally we demonstrate the effectiveness of our approach across a variety of continuous control tasks in Section 5 before concluding in Section 6, where we also mention some exciting future work.

2 Related Work

The Dyna algorithm [50] is a canonical approach for model based reinforcement learning (MBRL), based on the principle of ‘trying things in your head’, using an internal model of the world. In its original form, Dyna contained an exploration bonus for each state-action pair, proportional to this uncertainty measure [49]. A decade later, principled approaches were proposed for exploration [3], yet their guarantees were only possible in discrete settings with a finite number of states.

In recent times there has been great deal of progress in MBRL, with success in Atari games [30, 44], and dexterous manipulation [36], while progress has been made on continuous control benchmarks [31, 9, 10, 28]. We note that our approach is orthogonal to these in the following ways: 1) our methods, or some subset of the methods we introduce, can be incorporated into existing MBRL algorithms; 2) we adhere to a strict Dyna style framework, and our methods are aimed at incorporating active learning into this paradigm. Other recent work [18] shows that Dyna can be extended to latent-state dynamics models. We note that our framework could be extended to this setting, however we instead focus on efficient data acquisition through active learning.

Active approaches (i.e., acquiring data in a principled manner) in MBRL have been considered previously, for example in [47, 40, 23]. Usually, ensembles of models are maintained, and an intrinsic reward, defined as some difference measure (i.e., KL-divergence, total variation) across the output of different models in the ensemble drives exploration. Such exploration might be ineffective however, as the policy may visit regions of the state space which have no relation to solving the task. This may also lead to unsafe exploration if deployed on a real robot. [23] bears similarities to our work in that it aims to improve model generalization through exploration, and has a criteria to trade-off exploration and exploitation. However their approach to exploration is purely novelty-seeking, and they collect data until they discover a model (or subset of models) that can fully model the MDP in question (i.e., when novelty falls below a predefined value). Once a model is discovered, they implement a policy

(through search) which exploits it to maximize performance. This has drawbacks mentioned above concerning wasted exploration in task-irrelevant states as well as unsafe exploration.

Also similar to our work is [1], who explicitly use an active learning approach for solving an RL task; a robot arm hitting a ball. Our work differs in three significant ways: 1) our policies are trained inside a model, not on the true environment; 2) we actively limit the number of trajectories per collection phase based on the data (they introduce a heuristic); 3) we do not have access to an off-line measure of generalization, and therefore must introduce an online-learning mechanism to seek out the optimal setting of the exploration parameter.

Approaches that produce Bayes-optimal exploration and exploitation with a model [13, 41] are also of relevance, however these methodologies do not scale well to high dimensional tasks [51].

Efficient exploration in environments with very sparse rewards also represents a relevant area of research. In such settings an agent follows its *curiosity*, quantified by either: 1) rewarding areas of the state-space that reduce uncertainty in some internal model (i.e., inverse or forward dynamics models) [8, 35, 25, 39, 5]; 2) rewarding un-visited areas of the state-space [3, 32, 37]. Our approach to exploration leverages a model ensemble, and sits in the former category.

There has also been a great deal of work on using maximum entropy principles as a means for exploration in model-free RL [19]. The aim is to find rewarding behavior whilst maximizing some measure of entropy. We differ from these works in both what entropy is maximized (action entropy v.s. model prediction entropy) and by not having task-specific, fixed temperature parameter that trades off reward and entropy/surprise. In later maximum entropy work, temperature selection has been formulated as a constrained optimization problem, such that performance is maximized subject to some minimum level of policy entropy [20]. In contrast, we select this parameter in an online manner that optimally improves our internal models.

3 Background

3.1 RL Policies & Markov Decision Processes

A Markov Decision Process (MDP, [2]) is a tuple $(\mathcal{S}, \mathcal{A}, P, R)$. Here \mathcal{S} and \mathcal{A} stand for the sets of states and actions respectively, such that for $s_t, s_{t+1} \in \mathcal{S}$ and $a_t \in \mathcal{A}$: $P(s_{t+1}|s_t, a_t)$ is the probability that the system/agent transitions from s_t to s_{t+1} given action a_t and $R(a_t, s_t, s_{t+1})$ is a reward obtained by an agent transitioning from s_t to s_{t+1} via a_t .

A policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ is a (possibly randomized) mapping (parameterized by $\theta \in \mathbb{R}^d$, e.g. weight of the neural network) from \mathcal{S} to \mathcal{A} . Policy learning is the task of optimizing parameters θ of π_θ such that an agent applying it in the environment given by a fixed MDP maximizes total (expected/discouted) reward over given horizon H . In this paper we consider MDPs with finite horizons.

In most practical applications the MDP is not known to the learner. In MBRL, we seek to use a dataset $\mathcal{D} = \{(s_t, a_t), s_{t+1}\}_{t=1}^N$ of observed transitions to train a dynamics/world model \hat{f}_ϕ parameterized by ϕ to approximate the true dynamics function $f(s_{t+1}|s_t, a_t)$ such that $\hat{f}_\phi(s_{t+1}|s_t, a_t) \approx f(s_{t+1}|s_t, a_t)$.

We aim to construct rich \mathcal{D} s for learning accurate enough models \hat{f}_θ , but only in those regions that are critical for training performant policies.

3.2 Sequential Model Based Optimization

Consider a black box function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ over some domain \mathcal{X} , whereby the goal is to find $x^* \in \mathcal{X}$ such that

$$x^* = \arg \max_{x \in \mathcal{X}} F(x) \quad (1)$$

Sequential Model Based Optimization (SMBO, [26]) is a model-based black box optimization method which seeks to learn a surrogate model \hat{F} , within the true model F . Using the surrogate model, it is possible to determine which data should be collected to discover the optimum point of the real black box function F . The surrogate model is sequentially updated with the data collected in order to obtain better estimates of the true F , and this process is repeated until convergence or limited to a set number of iterations.

Many MBRL algorithms follow this regime, by trying to model a true black box function F using a world model \hat{f}_ϕ parameterized by ϕ as follows:

$$F(\theta) \cong \sum_{t=0}^H \gamma_t R(s_t, \pi_\theta(s_t), \hat{f}_\phi(\pi_\theta(s_t), s_t)) \quad (2)$$

where $\gamma_t \in [0, 1]$ is a discount factor, actions are $\pi_\theta(s_t) = a_t$, and next states are generated by $s_{t+1} = \hat{f}_\phi(\pi_\theta(s_t), s_t)$. Subsequently, they seek to find a policy π parameterized by θ that maximizes the acquisition function, which is just $\max F(\theta)$, then they will generate rollouts in the environment by using this policy and adding noise to the actions, such as via a stochastic policy.

Taking this view, essentially the majority of existing MBRL algorithms are conducting a form of SMBO with a greedy acquisition function. Any exploration conducted is simply a post hoc addition of noise in the hope of injecting sufficient stochasticity, but there is no principled mechanism by which the model may escape local minima.

Our approach is markedly different. We seek to conduct *Active Learning*, whereby the goal is not just to find an optimal solution to F , but to learn an optimal surrogate model \hat{F}^* through the world model \hat{f}_{ϕ^*} . To do this, we train policies to maximize an acquisition function that trades-off reward and information, rather than just greedily maximizing reward. Given how the acquisition function in SMBO is an important component for finding optimal performance [29], it therefore makes sense to carefully design how we acquire data in the MBRL setting.

3.3 Active Learning

Active Learning considers the problem of choosing new data $\mathcal{D}' = \{x'_i, y'_i\}_{i=1}^m$, which is a subset of a larger dataset i.e., $\mathcal{D}' \in \mathcal{D} = \{x_i, y_i\}_{i=1}^n$, such that a model \mathcal{M} is most improved.

In traditional supervised learning, we usually have access to the entirety of \mathcal{D} for training, but in the active learning setting we only have access to its feature vectors $\{x_i\}_{i=1}^n$, and need to query an oracle to obtain the corresponding labels $\{y_i\}_{i=1}^n$ which incurs a cost. Active learning aims to reduce this cost by iterating through the set of unlabeled feature vectors $\{x_i\}_{i=1}^n$, and determining which subset $\{x'_i\}_{i=1}^m$ would produce the greatest improvement in the performance of \mathcal{M} should we train our models on the subset $\mathcal{D}' = \{x'_i, y'_i\}_{i=1}^m$. The ultimate goal is to achieve the highest performance with the fewest number of queries (i.e., at the lowest cost).

In reinforcement learning, there is a key difference to the standard active learning setting; we do not have direct access to a dataset \mathcal{D} , and must instead generate \mathcal{D} through placing a ‘sampling’ policy π_s in the environment, producing trajectories. We then assess these trajectories generated by π_s and determine how to best train this policy to obtain trajectories in some oracle/true environment that will benefit performance. Framing the problem this way is done in [1], where it is cheap to generate and evaluate trajectories, but expensive to obtain labels.

In the Dyna-style approach [50] that we adopt, this means training our sampling policy π_s in the world model exclusively, then using this policy to collect samples in the real world. Through the lens of active learning, it is therefore important to train π_s in our world model such that the resultant trajectories maximize our world model performance, therefore maximizing our final policy performance. This is because it is free, from a data collection perspective, to train in the world model, but acquiring data in the real environment incurs an expense.

Concretely, we wish to train a policy that performs well in the real environment in as few samples as possible, and we identify that having a robust and generalizable world model as being paramount to achieving this. We observe that state-action pairs which cause high uncertainty/disagreement in the world model are likely to be parts of the state space that our world model is poor at modeling (as observed in [17], where solving ‘difficult’ environment dynamics is important for learning optimal policies).

However we cannot simply target regions of disagreement; we are primarily concerned with maximizing policy performance, so wish to explore regions of disagreement that are also critical to solving the task. This becomes a classic explore/exploit dilemma; how much do we want our sampling policy to explore (i.e., find states with poor generalization) or exploit (i.e., visit known high value states) when acquiring new real samples. In order to manage this trade-off, we leverage an online active learning

approach similar to [38], where model generalization feedback from the gathered trajectories can be used to determine the degree to which we explore or exploit in the future.

3.4 Online Learning

Online Learning is a family of methods that is used when a learner tackles a decision-making task by taking actions $a \in \mathcal{A}$, whilst learning from a sequence of data z_1, z_2, \dots, z_N that arrive in incremental rounds n . The learner must take an action at the beginning of each round, and the aim is to minimize cumulative regret, which is usually defined as the difference in some loss $\ell(a, z_n)$ between the actual action taken and the optimal action that could have been taken at that round:

$$\sum_{n=0}^N \left(\ell(a_n, z_n) - \min_{a \in \mathcal{A}} \ell(a, z_n) \right).$$

At each round n the learner does not know what the consequence of some action a will be, and only receives feedback after submitting its chosen action a_n . Based on this feedback the learner then updates how it selects actions in the future. In our approach we consider the set of actions \mathcal{A} to be the degree to which our policy explores, and the loss ℓ to be a normalized generalization error from the data collected. The aim is to therefore ensure that generalization error (i.e., RMSE on the new data) is maximized at each round. Since the task of maximizing generalization is both noisy and stochastic (i.e., the optimal degree of exploration may vary as we collect data), careful design of this algorithm is required.

4 Ready Policy One

Here we introduce our main algorithm, *Ready Policy One* (RP1). The key differences between RP1 and existing state of the art MBRL methods are as follows:

1. By taking an Active Learning approach rather than focusing on greedy optimization, RP1 seeks to directly learn the best *model*, rather than learning the best *policy*, and indirectly learning the best model to achieve this objective.
2. We introduce a principled Online Learning-inspired framework, allowing RP1 to adapt the level of exploration in order to optimally improve the model in the fewest number of samples possible.
3. We introduce a mechanism to stop gathering new samples in any given collection phase when the incoming data resembles what we have already acquired during that phase.

The algorithm begins in a similar fashion to other MBRL methods, by sampling initial transitions with a randomly initialized policy. In the Dyna framework, a policy is then trained inside \hat{f}_ϕ , and then subsequently used to gather new data. Typically, random noise is added to the policy to induce exploration. Other methods consider a hybrid objective. In RP1, consider training a *sampling policy*, parameterized by θ_t , to optimize the following objective:

$$\pi_{\theta_t} = \max[\mathbb{E}_{\tau \sim \pi_{\theta_t}} [(1 - \lambda)\mathcal{R}(\tau) + \lambda\sigma(\mathcal{R}(\tau))]] \quad (3)$$

where $\mathcal{R}(\tau) = \sum_{i=0}^H r_i$, $r_i = R(s_i, a_i, s_{i+1})$ and $\sigma(\mathcal{R}(\tau)) = \sum_{i=1}^H \sqrt{\frac{\sum_{j=1}^M (r_i^j - \bar{r}_i)^2}{M-1}}$. This λ value is chosen before training the policy in the model, and is selected using an online learning algorithm mechanism detailed in 4.2. Hence λ defines the relative weighting of reward and reward variance, with $\lambda = 0$ training a policy that only maximizes expected return (model-guided exploitation)*, and $\lambda = 1$ training a policy that only maximizes variance/disagreement per time step (model-guided exploration). In reality we limit λ to $[0.0, 0.5]$ as we wish any exploration to be guided by the reward signal. As a consequence, there is no need to train a separate ‘exploit’ policy, since we find policies trained in this way provide significant gains over commonly used approaches. This is mirrors the fact that MaxEnt strategies obtain high-performance in deterministic environments [15].

* $\lambda = 0$ corresponds to the same objective as in prior MBRL work, such as [31].

4.1 Information Gain in the Model as Maximum Entropy

Our objective in Equation 3 is inspired by curiosity driven exploration via model disagreement [40]. When shown a new unlabeled feature x' comprising a state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$ pair, the model has an existing ‘prior’ $p(r|x', \mathcal{D})$ where r is the predicted reward and \mathcal{D} is the data seen so far. After obtaining the label s' we have a new datapoint D' (i.e., $D' = (s, a, s')$), and can update the world model, producing a ‘posterior’ $p(r|s, a, \{\mathcal{D} \cup D'\})$. In our case, reward is a deterministic function of this triple (i.e., $r = R(D')$, see Appendix B for details). As such, we define the *Information Gain* (IG) in the reward as the KL divergence between the model posterior after observing D' and its respective prior at x' , as follows:

$$\mathcal{IG}(r; x') = D_{\text{KL}}[p(r|x', \{\mathcal{D} \cup D'\}) || p(r|x', \mathcal{D})] \quad (4)$$

$$= \int p(r|x', \{\mathcal{D} \cup D'\}) \log \frac{p(r|x', \{\mathcal{D} \cup D'\})}{p(r|x', \mathcal{D})} dr. \quad (5)$$

We observe that to maximize $\mathcal{IG}(r; x')$, we must sample x' appropriately, which is our only degree of freedom (through policy training). Because we cannot tractably calculate this quantity for all D' , one approach to maximize information gain is to ensure that the prior assigns low mass to all regions of r (i.e., minimize the denominator over all r). In order to do this we would select the improper prior over the continuous variable r^* . In our setting however, the model takes the form of an empirical Gaussian distribution, formed from the sample mean and variance of the individual models in the ensemble. Therefore we would like to like $p(r|x', \mathcal{D})$ such that the following is minimized:

$$D_{\text{KL}}[p(r|x', \mathcal{D}) | p_0(r)] \quad (6)$$

where $p_0(r)$ is an improper prior. The only way that Equation 6 is minimized is when the differential entropy of $p(r|x', \mathcal{D})$ is maximized. The differential entropy of a Gaussian is well known, given by $h(x) = \ln(\sqrt{2\pi}\sigma) + \frac{1}{2}$. Therefore to maximize the entropy of $p(r|x', \mathcal{D})$, and maximize information gain in the face of uncertainty, we need to maximize its variance. This is achieved by training policies that generate trajectory tuples $x' = (s, a)$ which cause high variance in the model, and is analogous to information based acquisition functions in active learning [33].

This motivates the second term of Equation 3, where we show the objective function for our exploration policies. Essentially, we are seeking to maximize information gain *in the model* through maximizing model entropy over the reward. This is in contrast to other maximum entropy approaches, which seek to maximize entropy over the action distribution [19], aiming to succeed at the task while acting as randomly as possible.

It is also possible to maximize information gain for next state predictions (as opposed to rewards), and this is similar to the approach in [40]. However in practice we find that maximizing reward variance results in better performance (see Section 5).

4.2 Online Learning Mechanism

We use the Exponential Weights framework to derive an adaptive algorithm for the selection of λ . In this setup we consider k experts making recommendations at the beginning of each round. After sampling a decision $i_t \in \{1, \dots, k\}$ from a distribution $\mathbf{p}^t \in \Delta_k$ with the form $\mathbf{p}^t(i) \propto \exp(\ell_t(i))$ the learner experiences a loss $l_{i_t}^t \in \mathbb{R}$. The distribution \mathbf{p}^t is updated by updating ℓ_t as follows:

$$\ell_{t+1}(i) = \begin{cases} \ell_t(i) + \eta \frac{l_{i_t}^t}{\mathbf{p}^t(i)} & \text{if } i = i_t \\ \ell_t(i) & \text{o.w.} \end{cases} \quad (7)$$

For some step size parameter η .

In our case we consider the case when the selection of λ_t is thought as choosing among k experts which we identify as the different values $\{\lambda_i\}_{i=1}^k$. The loss we consider is of the form $l_{i_t}^t = \hat{G}_{\phi_t}(\theta_{t+1})$, where $G_{\phi_t}(\theta_{t+1})$ is the RMSE of the model under parameters ϕ_t on data collected using $\pi_{\theta_{t+1}}$, and θ_{t+1} is the parameter of the policy trained under the choice λ_{i_t} , after incorporating into the model the data collected using the previous policy π_{θ_t} . We then perform a normalization of G (see

*This follows the proposal in [34] when selecting priors in the face of uncertainty.

Appendix B for details), hence \hat{G} . Henceforth we denote by \mathbf{p}_λ^t the exponential weights distribution over λ values at time t .

Algorithm 1 Online Learning Mechanism

Input: step size η , number of timesteps T .

Initialize: \mathbf{p}_λ^1 as a uniform distribution.

for $t = 1, \dots, T - 1$ **do**

1. Select $i_t \sim \mathbf{p}_\lambda^t$ and $\lambda_t = \lambda_{i_t}$.

2. Use Equation 7 to update \mathbf{p}_λ^t with

$$l_{i_t}^t = \hat{G}_{\phi_t}(\theta_{t+1})$$

Our version of Exponential Weights algorithm also known as Exponentially Weighted Average Forecaster [6] is outlined in Algorithm 1.

In practice and in order to promote more effective exploration over λ values we sample from a mixture distribution where \mathbf{p}_λ^t is not proportional to $\exp(\ell_t)$ but it is a mixture between this exponential weights distribution and the uniform over $[k]$. In other words, let $\epsilon > 0$ be a small parameter. With probability $1 - \epsilon$ the produce i_t as a sample from an exponential weights distribution proportional to $\exp(\ell_t)$, and with probability ϵ it equals a uniform index from $1, \dots, k$.

4.3 Diverse Sample Collection

Consider the problem of online data acquisition from a policy in an environment. At each timestep we receive a set of datapoints $\{x_1, \dots, x_H\} \sim \pi_\theta$ corresponding to the concatenation of each state and action in a trajectory. At timestep t we have a dataset $X_t = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$, where $X_t \in \mathbb{R}^{d \times n}$ sampled from the sampling policy. We represent this data in the form of the Singular Value Decomposition (SVD) of the symmetric matrix, $\text{Cov}_t = \frac{1}{n} X_t X_t^\top = \mathbf{Q}_t^\top \Sigma_t \mathbf{Q}_t \in \mathbb{R}^d$.

Equipped with this representation, we take the top k eigenvalues λ_i of Cov_t , where k is smallest such that: $\sum_{i=1}^k \lambda_i \geq (1 - \delta) \sum_{i=1}^d \lambda_i$ for some parameter $\delta > 0$, and take $n_t = k$. Next we take the corresponding eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^d$ and let $\mathbf{U} \in \mathbb{R}^{d \times k}$ be obtained by stacking them together. We define the Active Subspace [12] $\mathbf{U}^{\text{act}} \in \mathbb{R}^{d \times k}$ as $\mathcal{L}_{\text{active}} \stackrel{\text{def}}{=} \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$. \mathbf{U}^{act} is an orthonormal basis of $\mathcal{L}_{\text{active}}$.

We use \mathbf{U}^{act} to evaluate new data. After we collect n' new samples $V_{t+1} \in \mathbb{R}^{d \times n'}$, we form a covariance matrix with this new data as $\frac{1}{n'} V_{t+1} V_{t+1}^\top \in \mathbb{R}^{d \times d}$ and project it onto \mathbf{U}^{act} . We define the residual at timestep t , r_t , as follows:

$$r_t = \frac{\text{tr}(V_{t+1} V_{t+1}^\top - \mathbf{U} \mathbf{U}^\top V_{t+1} V_{t+1}^\top \mathbf{U} \mathbf{U}^\top)}{\text{tr}(V_{t+1} V_{t+1}^\top)} \quad (8)$$

Where tr denotes the trace operator. After evaluating r_t we append the new data V_{t+1} to X_t to form $X_{t+1} \in \mathbb{R}^{d \times (n+n')}$. Intuitively, r_t tells us how much of the new data could not be explained by the principal components in the data collected thus far. We stop collecting data and proceed to retrain the model once $r_t < \alpha$, where α is a proxy for δ . The full procedure is presented in Algorithm 2.

Let q_t be the probability that at timestep t , step 4. of Algorithm 2 is executed (i.e., $q_t = \text{P}(r_t < \alpha)$). The evolution of q_t operates in roughly two phases. First, the algorithm tries to collect data to form an accurate estimate of the covariance matrix Cov_t and a stable estimator \mathbf{U}^{act} . During this phase, q_t is small as it is roughly the probability of two random samples $X, X' \in \mathbb{R}^{d \times n'}$ aligning. After t_0 steps when the algorithm has built a stable estimator of \mathbf{U}^{act} , the stopping probability stabilizes to a value q^* that solely depends on the trajectories intrinsic noise. Both the length of t_0 and the magnitude of q^* scale with the trajectories' noise. If the trajectories have little noise both t_0 and q^* are small. On the other hand, if the trajectories have high noise, the early stopping mechanism will take longer to trigger.

Algorithm 2 Early Stopping Mechanism

Input: thresholds α, δ , maximum number of samples T .

Initialize: training set $X = \emptyset$

Collect initial samples $X_1 = \{x_1, \dots, x_H\}$.

for $t = 2, \dots, T - 1$ **do**

1. Compute Active Subspace \mathbf{U}^{act} as the result of stacking together some orthonormal basis of $\mathcal{L}_{\text{active}} \stackrel{\text{def}}{=} \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ where the vectors \mathbf{u}_i correspond to the top k eigenvalues of the covariance matrix Cov_t .
 2. Produce samples V_{t+1} via the sampling policy
 3. Calculate the residual r_t using Equation 8.
 4. Stop collecting data if $r_t < \alpha$
-

This dynamic approach to determining the effective ‘batch size’ of the incoming labeled data is similar to [7], whereby feedback from unsupervised learning is used to control the amount of data collected per batch. However, we do this in a more instantaneous fashion, leveraging data collected so far to determine when to stop.

4.4 The Algorithm

Algorithm 3 RP1: Ready Policy One

Input: Number of initial samples N_0 , number of ongoing samples N_t , number of policies in the ensemble M , number of time steps T .

Initialize: Initial World Model f_0 comprised of M models.

Collect N_0 samples with a random policy and initialize data set $\mathcal{D}_0 = \{(s_t, a_t), s_{t+1}\}_{t=1}^{N_0}$.

for $t = 1, \dots, T - 1$ **do**

1. Train \hat{f}_{t-1} with \mathcal{D}_t to derive \hat{f}_t .
 2. Select λ using Algorithm 1.
 3. Train exploration policy π_{ϕ_t} using Equation 3 in \hat{f}_t .
 4. Collect new samples $\mathcal{D}_{\text{new}} = \{(s_t, a_t), s_{t+1}\}_{t=1}^{N_t}$ in the environment with π_{ϕ_t} , where N_t is defined as the number of time steps required for Algorithm 2 to return.
 5. $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \mathcal{D}_{\text{new}}$
 6. Set $\alpha_t = \mathcal{L}(\hat{f}_t(\mathcal{D}_{\text{new}}))$
-

We now present our algorithm: Ready Policy One (RP1). At each iteration we select a policy objective (λ) to maximize sample utility, and train a policy on this objective. The policies are trained using the original PPO [46] loss function, but we use the training approach in [45] as this combination delivered more robust policy updates.*

Once the policy has been trained inside the model, we use it to generate samples in the real environment. These samples continue until our early stopping mechanism is triggered, and we have sufficiently diverse data to retrain the model. The full procedure is outlined in Algorithm 3.

The overall aim is to therefore determine which part of the model space is both high value and unknown, so that our trained sampling policy can obtain enough data samples pertaining to those regions of the environment.

5 Experiments

The primary goal of our experiments is to evaluate whether our active learning approach for MBRL is more sample efficient than existing approaches. In particular, we test RP1 on a variety of continuous control tasks from the OpenAI Gym [4], namely: HalfCheetah, Ant, Swimmer and Hopper, which are commonly used to test MBRL algorithms. For specifics, see Appendix B. In order to produce robust results, we run all experiments for ten random seeds, more than typically used for similar analyses [24].

*Full implementation details can be found in Appendix B.

Rather than individual algorithms, we compare against the two approaches most commonly used in MBRL:

- **Greedy:** We train the policy to maximize reward in the model, and subsequently add noise to discover previously unseen states. This is the approach used in ME-TRPO [31].
- **Variance + Reward (V+R):** We train the policy with $\lambda = 0.5$, producing a fixed degree of priority for reward and model entropy. This resembles methods with exploration bonuses such as [25].

We note that these baselines are non-trivial. In particular, ME-TRPO is competitive with state of the art in MBRL. In fact, for two of the tasks considered (Swimmer and Hopper) it outperformed all other approaches in a recent paper benchmarking MRBL methods [51].

We also compare against the same policy gradients algorithm as a model free baseline, which we train for 10^6 , 5×10^6 and 10^7 timesteps. This provides an indication of the asymptotic performance of our policy, if trained in the true environment.

Table 1: Median best performance at a given timestep for ten seeds. Bold indicates the best performing algorithm. T1 corresponds to the t-stat for RP1 vs. Greedy, T2 corresponds to the t-stat for RP1 vs. V+R. * indicates $p < 0.05$.

	Timesteps	Greedy	V+R	RP1	T1	T2
HalfCheetah	10^4	-0.95	-1.1	100.51	5.39*	4.34*
Ant	10^4	94.72	95.07	113.63	4.02*	3.08*
Swimmer	10^4	1.08	1.07	3.24	1.2	1.98
Hopper	10^4	76.5	139.03	322.22	4.32*	3.42*
HalfCheetah	10^5	260.92	283.27	390.49	3.89*	2.62*
Ant	10^5	186.36	217.08	238.4	3.83*	3.11*
Swimmer	10^5	61.76	62.89	64.19	-0.11	-1.09
Hopper	10^5	487.25	570.09	619.73	3.52*	2.21

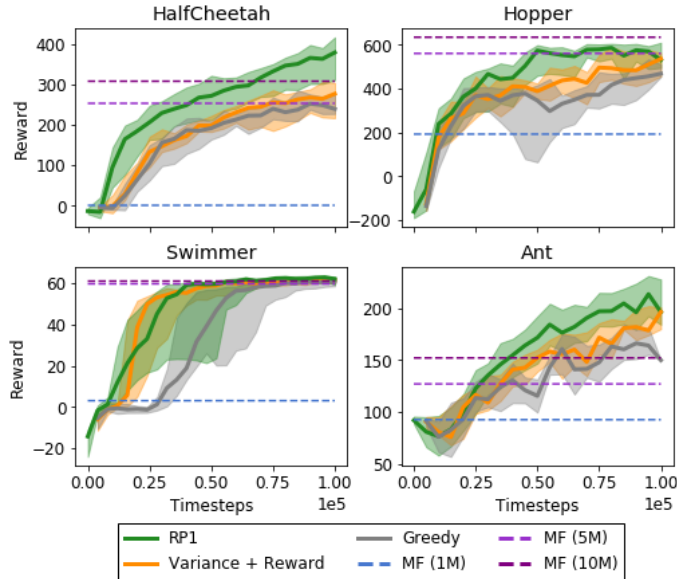


Figure 1: Median performance across 10 seeds. Shaded regions correspond to the Inter-Quartile Range.

Table 1 and Fig 1 show the main results, where RP1 outperforms both the greedy baseline and the fixed variance maximizing (V+R) approach. Furthermore, we perform Welch’s unequal variances t-test and see that in most cases the results are statistically significant, aside from Swimmer at 10^5 timesteps where all three methods have converged to the optimal solution. In addition, we observe that RP1 is able to achieve strong performance vs. model-free in fewer timesteps than the existing baselines.

Interestingly, we see that simply adding a fixed entropy term (V+R) into the reward function gives improved performance over the baseline greedy approach. This corroborates with findings in [17], where there is, in most tasks, a correlation between regions that are difficult to predict and regions of high reward. However our findings also suggest that this is not always the case, and having the ability to adjust how much we focus on such regions of disagreement is vital. We hypothesize that the fixed V+R approach may collect too many high-variance samples for certain tasks, since we do not tune λ , nor limit batch size. As a result, the trajectories gathered do not necessarily result in strong policy performance, unlike RP1, which aims to maximize policy performance through the data collected.

We support this hypothesis with Fig. 2 where, we show the normalized change in reward for different λ values in each task. In particular, we observe for HalfCheetah, Swimmer and Ant, a greater focus on uncertainty appears positively related to faster learning. However, the opposite is true for Hopper. The benefit of our mechanism is the ability to dynamically learn this preference, and thus adapt to the current optimization landscape.



Figure 2: Mean one-step policy improvement after a given λ value, for all ten seeds of RP1.

Next we study the choice of model variance used in the reward function. Other work, such as [40] use the variance over the next state prediction, whereas RP1 uses the variance over the reward. Fig 3 and Table 2 show that next state variance is a less effective approach. This is likely due to over-emphasis of regions of the state space that are inherently hard to predict, but do not impact the ability to solve the task [43].

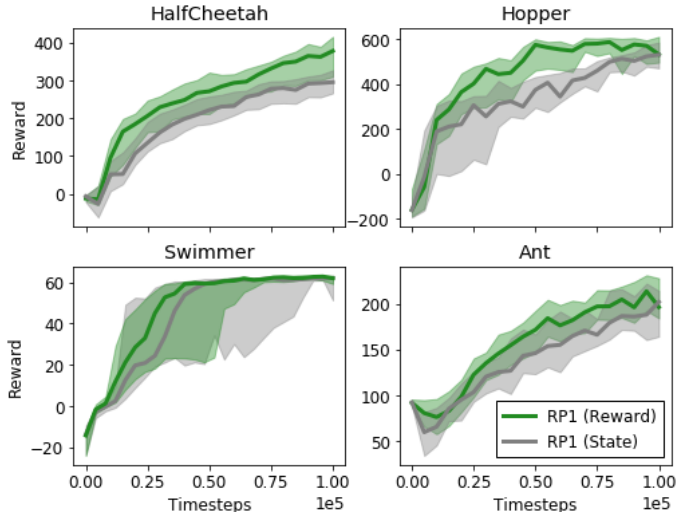


Figure 3: Median performance across 10 seeds. Shaded regions correspond to the Inter-Quartile Range.

Table 2: Study for the choice of error to maximize. Results show the median best performance at 10^5 timesteps for ten seeds. The highest performing value for each environment is bolded.

	HalfCheetah	Ant	Swimmer	Hopper
State	319.41	214.33	63.47	549.19
Reward	390.49	238.4	64.19	619.73

Finally, we consider the individual components of the RP1 algorithm. We evaluate two variants: RP1 ($\lambda = 0$), where we remove the online learning mechanism and train a greedy policy, and RP1 (No EarlyStop), where we remove the early stopping mechanism and use a fixed batch size. Results are shown in Table 3, and Figs 4 and 5 in Appendix A.

Table 3: Ablation study for the key components of RP1. Results show the median performance at 10^5 timesteps for 10 seeds. The highest performing value for each environment is bolded.

	HalfCheetah	Ant	Swimmer	Hopper
Baseline	283.27	186.36	62.89	487.25
RP1 ($\lambda = 0$)	319.14	223.21	41.14	603.52
RP1 (No EarlyStop)	247.82	197.95	41.04	595.77
RP1	390.49	238.4	64.19	619.73

We observe that the improvements attributed to RP1 are not down to any single design choice, and the individual components complement each other to provide significant overall gains. For example, by conducting purely noise-based exploration ($\lambda = 0$), we lose the flexibility to target specific regions of the state-space. On the other hand, by removing our early stopping mechanism (No EarlyStop), we acquire a trajectory dataset for our model that has too much redundant data, reducing sample efficiency. Nonetheless, we believe adding either of these components to existing MBRL methods, which either have a fixed temperature parameter (λ) or fixed data collection batch size, would lead to performance gains.

6 Conclusion and Future Work

We presented Ready Policy One (RP1), a new approach for Model-Based Reinforcement Learning (MBRL). RP1 casts data collection in MBRL as an active learning problem, and subsequently seeks to acquire the most informative data via an exploration policy. Leveraging online learning techniques, the objective function for this policy adapts during optimization, allowing RP1 to vary its focus on the often fruitful reward function. We showed in a variety of experiments that RP1 significantly increases sample efficiency in MBRL, and we believe it can lead to new state of the art when combined with the latest architectures.

We are particularly excited by the many future directions from this work. Most obviously, since our method is orthogonal to other recent advances in MBRL, RP1 could be combined with state of the art probabilistic architectures [9], or variational autoencoder based models [18, 22].

In addition, we could take a hierarchical approach, by ensuring our exploration policies maintain core behaviors but maximize entropy in some distant unexplored region. This would require behavioral representations, and some notion of distance in behavioral space [11], and may lead to increased sample efficiency as we could better target specific state action pairs.

References

- [1] T. Akiyama, H. Hachiya, and M. Sugiyama. Efficient exploration through active learning for value function approximation in reinforcement learning. *Neural Networks*, 23(5):639 – 648, 2010.
- [2] R. Bellman. A Markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
- [3] R. I. Brafman and M. Tennenholtz. R-max - a general polynomial time algorithm for near-optimal reinforcement learning. *J. Mach. Learn. Res.*, 3(null):213231, Mar. 2003.
- [4] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI Gym, 2016.
- [5] Y. Burda, H. Edwards, D. Pathak, A. J. Storkey, T. Darrell, and A. A. Efros. Large-scale study of curiosity-driven learning. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- [6] N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.

- [7] S. Chakraborty, V. Balasubramanian, and S. Panchanathan. Dynamic batch mode active learning. In *2011 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011*, pages 2649–2656, 2011.
- [8] N. Chentanez, A. G. Barto, and S. P. Singh. Intrinsically motivated reinforcement learning. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 1281–1288. MIT Press, 2005.
- [9] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 4754–4765. 2018.
- [10] I. Clavera, J. Rothfuss, J. Schulman, Y. Fujita, T. Asfour, and P. Abbeel. Model-based reinforcement learning via meta-policy optimization. In *2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings*, pages 617–629, 2018.
- [11] J. D. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine. Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. In *ICML*, 2018.
- [12] P. G. Constantine, E. Dow, and Q. Wang. Active subspace methods in theory and practice: Applications to kriging surfaces. *SIAM Journal on Scientific Computing*, 36(4):A1500–A1524, 2014.
- [13] M. P. Deisenroth and C. E. Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on International Conference on Machine Learning, ICML11*, page 465472, Madison, WI, USA, 2011. Omnipress.
- [14] Y. Duan, X. Chen, R. Houthoof, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement learning for continuous control. *CoRR*, abs/1604.06778, 2016.
- [15] B. Eysenbach and S. Levine. If maxent rl is the answer, what is the question?, 2019.
- [16] V. Feinberg, A. Wan, I. Stoica, M. I. Jordan, J. E. Gonzalez, and S. Levine. Model-based value estimation for efficient model-free reinforcement learning. *CoRR*, abs/1803.00101, 2018.
- [17] D. Freeman, D. Ha, and L. Metz. Learning to predict without looking ahead: World models without forward prediction. In *Advances in Neural Information Processing Systems 32*, pages 5380–5391. Curran Associates, Inc., 2019.
- [18] D. Ha and J. Schmidhuber. Recurrent world models facilitate policy evolution. In *Proceedings of the 32Nd International Conference on Neural Information Processing Systems, NeurIPS’18*, pages 2455–2467, USA, 2018. Curran Associates Inc.
- [19] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1861–1870, Stockholm, Sweden, 2018.
- [20] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine. Soft actor-critic algorithms and applications. *CoRR*, abs/1812.05905, 2018.
- [21] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [22] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent dynamics for planning from pixels. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, pages 2555–2565, 2019.
- [23] M. Henaff. Explicit explore-exploit algorithms in continuous state spaces. In *Advances in Neural Information Processing Systems 32*, pages 9372–9382. Curran Associates, Inc., 2019.
- [24] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger. Deep reinforcement learning that matters. *CoRR*, abs/1709.06560, 2017.
- [25] R. Houthoof, X. Chen, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems 29*. 2016.

- [26] F. Hutter, H. H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In C. A. C. Coello, editor, *Learning and Intelligent Optimization*, pages 507–523, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [27] A. Ilyas, L. Engstrom, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Are deep policy gradient algorithms truly policy gradient algorithms? *CoRR*, abs/1811.02553, 2018.
- [28] M. Janner, J. Fu, M. Zhang, and S. Levine. When to trust your model: Model-based policy optimization. *CoRR*, abs/1906.08253, 2019.
- [29] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- [30] Ł. Kaiser, M. Babaeizadeh, P. Miłoś, B. Osiański, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, A. Mohiuddin, R. Sepassi, G. Tucker, and H. Michalewski. Model based reinforcement learning for atari. In *International Conference on Learning Representations*, 2020.
- [31] T. Kurutach, I. Clavera, Y. Duan, A. Tamar, and P. Abbeel. Model-ensemble trust-region policy optimization. In *International Conference on Learning Representations*, 2018.
- [32] M. Lopes, T. Lang, M. Toussaint, and P. yves Oudeyer. Exploration in model-based reinforcement learning by empirically estimating learning progress. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 206–214. Curran Associates, Inc., 2012.
- [33] D. J. C. MacKay. Information-based objective functions for active data selection. *Neural Computation*, 4(4):590–604, July 1992.
- [34] D. J. C. MacKay. *Information Theory, Inference & Learning Algorithms*. Cambridge University Press, USA, 2002.
- [35] S. Mohamed and D. J. Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS15*, page 21252133, Cambridge, MA, USA, 2015. MIT Press.
- [36] A. Nagabandi, K. Konoglie, S. Levine, and V. Kumar. Deep dynamics models for learning dexterous manipulation, 2019.
- [37] G. Ostrovski, M. G. Bellemare, A. van den Oord, and R. Munos. Count-based exploration with neural density models. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML17*, page 27212730. JMLR.org, 2017.
- [38] T. Osugi, D. Kun, and S. Scott. Balancing exploration and exploitation: A new algorithm for active machine learning. In *Proceedings of the Fifth IEEE International Conference on Data Mining, ICDM 05*, page 330337, USA, 2005. IEEE Computer Society.
- [39] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 2778–2787, 2017.
- [40] D. Pathak, D. Gandhi, and A. Gupta. Self-supervised exploration via disagreement. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 5062–5071, 2019.
- [41] N. Sajid, P. J. Ball, and K. J. Friston. Active inference: demystified and compared. *arXiv e-prints*, page arXiv:1909.10863, Sep 2019.
- [42] J. Schmidhuber. Curious model-building control systems. [*Proceedings*] *1991 IEEE International Joint Conference on Neural Networks*, pages 1458–1463 vol.2, 1991.
- [43] J. Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (19902010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, Sep. 2010.
- [44] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering Atari, Go, chess and shogi by planning with a learned model, 2019.
- [45] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International Conference on Machine Learning (ICML)*, 2015.

- [46] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms, 2017.
- [47] P. Shyam, W. Jaskowski, and F. Gomez. Model-based active exploration. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 5779–5788, 2019.
- [48] Y. Sun, F. Gomez, and J. Schmidhuber. Planning to be surprised: Optimal Bayesian exploration in dynamic environments. In *Proceedings of the 4th International Conference on Artificial General Intelligence, AGI11*, page 4151, Berlin, Heidelberg, 2011. Springer-Verlag.
- [49] R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *In Proceedings of the Seventh International Conference on Machine Learning*, pages 216–224. Morgan Kaufmann, 1990.
- [50] R. S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART Bull.*, 2(4):160163, July 1991.
- [51] T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, P. Abbeel, and J. Ba. Benchmarking model-based reinforcement learning. *CoRR*, abs/1907.02057, 2019.

A Additional Ablation Studies

In this section we show the full results from the ablation study in Table 3.

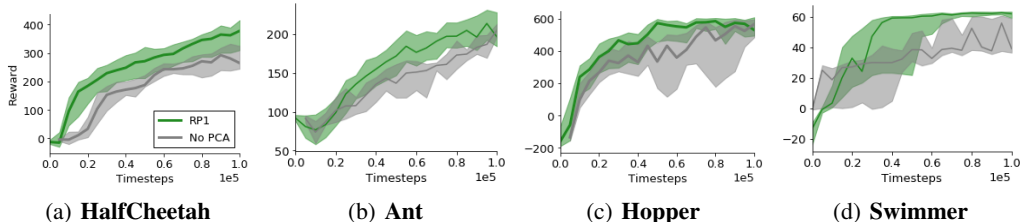


Figure 4: Ablation study where we consider removing the early stopping mechanism. All results show the median performance across ten seeds.

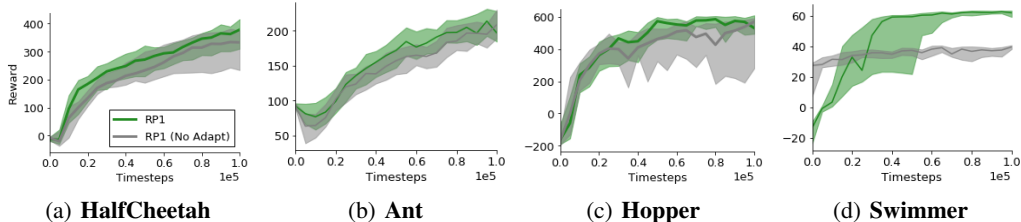


Figure 5: Ablation study where we consider removing the adaptive mechanism. All results show the median performance across ten seeds.

B Implementation Details

In terms of approach, we follow ME-TRPO [31] with some adjustments. Instead of using TRPO loss [45], we leverage the first-order approximation loss in PPO [46]. We do not apply parameter noise, and we modify policy training slightly by ensuring at least 10 updates are performed before termination is considered, which we find helps to improve convergence. We do not apply GAE nor the overall training approach in [46], as this introduced instabilities. We instead found that generating large batch sizes gave better and more consistent performance, which corroborates the findings in [27] with respect to true policy gradients.

We augment each environment with an additional state which contains velocity information. This is also done in the ‘FixedSwimmer’ environment in [51], and allows us to infer the reward from the states directly. It must also be noted that in the original ‘rllab’ [14] environments used in [31], one of the observable states was the velocity state used to calculate rewards, and we therefore mirror this in our OpenAI Gym [4] implementation; we do not anticipate there to be any problem integrating reward prediction with our framework. Furthermore, we provide this state in both the model-free and model-based benchmarks to ensure there is no advantage, and do not notice any noticeable improvement in the model-free setting when this is provided; we hypothesize that some close proxy to the true velocity state already exists in the original state-space. We remove contact information from all environments, and instead of ‘Swimmer-v2’ we use the aforementioned ‘FixedSwimmer’, since this can be solved by our policy in a model-free regime. We remove early stopping from Hopper since we found it was necessary for convergence, but left the early stopping in for Ant since it was possible to train performant policies. We train the policy for 100 time steps in HalfCheetah and Ant, and for 200 time steps in Swimmer and Hopper. In experiments without the early stopping mechanism, data collection defaults to 3,000 timesteps per iteration. Full hyperparameter values can be found in Table 5.

We use the following approach to normalize G to produce \hat{G} ; for convenience, we write $\hat{G}_{\phi_t}(\theta_{t+1})$ as \hat{G}_t .

$$\hat{G}_t = \frac{G_t - \frac{1}{5} \sum_{\tau=t-5}^{t-1} (G_\tau)}{l_{\text{val}}} \tag{9}$$

where l_{val} is the final model validation loss from the iteration $t - 1$.

Attention should be drawn to the α parameter used to determine early stopping in Algorithm 2. For the tasks we test on, we choose to fix this to 0.0005, and therefore do not tune it to be task specific. We found that at this setting of α , the early stopping mechanism generally collects significantly fewer than the default 3,000 samples (which acts as an upper bound in RP1), but can still expand the batch size collected to the full amount where appropriate (i.e., under policies that provide non-homogenous trajectories).

Table 4: Hyperparameters used in the Policy

Hyperparameter Name	Value
Optimizer	Adam
Learning Rate	3e-4
Loss	PPO
Discount Factor	0.99
Batch Size	50,000
Epochs per Batch	10
ϵ -clip	0.2
Default Action σ	0.5
Hidden Layer Size	32
Number of Hidden Layers	2
Activation Function	ReLU

Table 5: Hyperparameters used in the World Model

Hyperparameter Name	Value
Optimizer	Adam
Learning Rate	1e-3
Train/Validation Split	2:1
Number of Models	5
Batch Size	1,024
Hidden Layer Size	1,024
Number of Hidden Layers	2
Activation Function	ReLU
Early Stopping α in PCA	0.0005