

---

# Bayesian-Adaptive Deep Reinforcement Learning using Model Ensembles

---

**Gilwoo Lee**      **Jeongseok Lee**      **Brian Hou**      **Aditya Mandalika**  
gilwoo@cs.uw.edu    jslee02@cs.uw.edu    bhoul@cs.uw.edu    adityavk@cs.uw.edu

## 1 Introduction

Although model-free deep reinforcement learning algorithms have shown tremendous success in a wide range of tasks, such as simulated control problems in OpenAI Gym [1] and games like Go [2], they face fundamental challenges in their application to physical control problems on robotic hardware. The high sample complexity of current methods results in long training periods if applied directly to hardware, even for simple grasping tasks [3]. Furthermore, the safety of the robot or surrounding environment when gathering data from the real world may be difficult to guarantee while the policy is training.

Model-based reinforcement learning techniques offer an opportunity to overcome these issues by taking advantage of simulations of the real systems, which can generate training data faster than real time. The primary challenge with model-based techniques is the discrepancy between simulation and the real world. Simplified models, inaccurate or uncertain parameters that govern the dynamics of the model, and other unmodeled disturbances and noise can render the policies learned on the simulated model ineffective on the real system.

We propose a model-based algorithm that learns a universal policy for Bayes-Adaptive Markov Decision Process (BAMDP) that is robust and optimal to model uncertainty and disturbances. Although the real system’s exact physical parameters are unknown, we can maintain a belief distribution over those parameters and provide this belief as an additional input to the policy. This enables the learned policy to be bold when confident in its model and cautious when there is high model uncertainty.

## 2 Related Work

Our work is closely related to QMDP [4, 5] which is a Q-value approximation method for Partially Observable Markov Decision Processes (POMDPs). QMDP assumes a fully-observable MDP after one action by approximating the Q-value at the current belief state  $b(s)$  as a weighted average over the fully-observable MDP’s Q-values,  $Q_{\text{MDP}}(b, a) = \sum_s b(s)Q(s, a)$ . The algorithm requires access to the (approximate) Q-function for the MDP. QMDP performs surprisingly well for many problems, but determinizing the belief distribution after one timestep breaks Bayes-optimality.

The BAMDP formulation that we consider is also similar to the POMDP formulation used in POMDP-lite [6], which assumes that the hidden state variables remain constant or only change deterministically. In our case, the hidden state variables correspond to the physics parameters  $\phi$ . Chen et al. show that this formulation is “equivalent to a set of fully observable Markov decision processes indexed by a hidden parameter”, which, in our case, is a discretization of  $\phi$ .

Robustness to model uncertainty has also been addressed by EPOpt [7], which learns a policy that maximizes the worst-case performance across multiple sampled MDPs. However, as is common with min-max techniques, EPOpt can be quite conservative (and therefore far from optimal) in its policy, especially when the model’s possible parameter-space is large. Our work closes this gap in EPOpt by maintaining a belief over the training MDPs to balance between optimality and robustness.



Figure 1: Graphical model for a Bayes-Adaptive Markov Decision Process.

Another recent approach to model uncertainty is UP-OSI [8], which utilizes online system identification to estimate the physical parameters of the system. However, UP-OSI only uses the mean estimate of those parameters as an additional input to the policy. Without a notion of belief and uncertainty in its parameter estimates, UP-OSI is prone to aggressively executing policies without being conservative toward remaining model uncertainty.

TODO: mention [9] and other bayesian work.

### 3 Bayes-Adaptive Reinforcement Learning

We follow the BAMDP framework (Figure 1), which assumes that a latent variable  $\phi$  governs the transition function of the underlying Markov Decision Process [10–12]. A BAMDP is defined by a tuple  $\langle \mathcal{S}', \mathcal{A}, P, P_0, R \rangle$ , where

- $\mathcal{S}' = \mathcal{S} \times \Phi$  is the set of hyper-states (state  $s$ , latent variable  $\phi$ ),
- $\mathcal{A}$  is the set of actions,
- $P(s', \phi' | s, \phi, a)$  is the transition function between hyper-states, conditioned on action  $a$  being taken in hyper-state  $(s, \phi)$ ,
- $P_0 \in \mathcal{P}(\mathcal{S} \times \Phi)$  combines the initial distribution over hyper-states,
- $R(s, \phi, a)$  represents the reward obtained when action  $a$  is taken in hyper-state  $(s, \phi)$ .

In our setting of robot control, the latent variables  $\phi$  are physics parameters such as mass and length of different robot links (TODO: better word than links) (JS: robot bodies?).

Our goal is to find the Bayes-optimal policy for the following Bellman equation:

$$V^*(b, s) = \max_a \left\{ R(s, b, a) + \gamma \sum_{s'} P(s', b' | s, b, a) V^*(b', s') \right\} \quad (1)$$

where  $b(\phi)$  is the belief over the set of latent physics parameters  $\phi \in \Phi$ ,

$$\begin{aligned} R(s, b, a) &= \sum_{\phi \in \Phi} b(\phi) R(s, \phi, a) \\ P(s', b' | s, b, a) &= \sum_{\phi \in \Phi} b(\phi) P(s', b' | s, \phi, a) \\ &= \sum_{\phi \in \Phi} b(\phi) P(s' | s, \phi, a) \sum_{\phi' \in \Phi} P(\phi' | s, \phi, a). \end{aligned}$$

TODO(?): mention PAC-MDP?

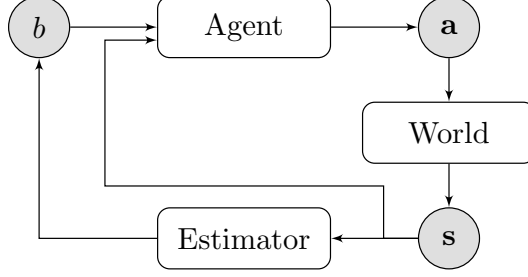


Figure 2: System structure

### 3.1 Belief Update

We make a simplification to the BAMDP formulation by assuming that the latent variable  $\phi$  is constant or changing deterministically according to a known transition function. As in POMDP-lite[6], we can then view the simplified BAMDP model as a set of MDPs that are indexed by the latent variable  $\phi$ .

The posterior probability of an MDP governed by latent physics variable  $\phi$  given a trajectory  $\tau_{0:H} = (s_0, a_0, \dots, s_H)$  is given by:

$$\begin{aligned}
 P(\phi_H | \tau_{0:H}) &= \frac{1}{Z} P(\phi_0) P(\tau_{0:H} | \phi_0) \\
 &= \frac{1}{Z} P(\phi_0) \prod_{t=0}^{H-1} P(s_{t+1}, \phi_{t+1} | s_t, \phi_t, a_t) \\
 &= \frac{1}{Z} P(\phi_0) \prod_{t=0}^{H-1} P(s_{t+1} | s_t, \phi_t, a_t) P(\phi_{t+1} | s_t, \phi_t, a_t)
 \end{aligned} \tag{2}$$

where  $Z$  is a normalizing constant and  $P(\phi_{t+1} | \cdot) = 1$  is the deterministic transition of  $\phi$ .

## 4 Bayes-Adaptive Policy Gradient

The main contribution of our work is a policy-gradient algorithm which produces a universal policy that maintains different strategies for different beliefs. The key idea is to augment the state with the belief and provide this augmented feature as input to the policy. The policy is a function of both state and belief:  $\pi : \mathcal{S} \times \mathcal{B} \rightarrow P(\mathcal{A})$ , where  $\mathcal{B}$  is the belief space. This allows the policy to differentiate between the beliefs while sharing globally effective strategies.

In order for this algorithm to work in practice, we need an effective representation of the belief space. If the space of MDPs  $\mathcal{M}$  is finite, the belief  $b$  can be a vector of size  $|\mathcal{M}|$ . If  $\mathcal{M}$  is infinite, we propose two methods to approximate the belief space:

- Sample  $K$  MDPs from the prior belief distribution  $b_0$  and use belief over these  $K$  MDPs
- Discretize the parameter space and use belief over the parameter space

In the first case, we use a vector of size  $K$  as  $b$ ;  $K$  should be large enough to approximate all MDPs with the  $K$  sampled ones. In the second case, assuming the latent parameters are independent, the belief is represented as a vector of size  $(N_1 + \dots + N_M)$  where  $M$  is the total number of parameters and  $N_m$  is the discretization of the  $m$ -th parameter. In this project, we have implemented the first one, and leave the second one for future work.

The algorithm is shown in Algorithm 1. At each iteration, we sample a number of MDPs from the prior distribution  $b_0$ , and sample a trajectory for each MDP using the current policy. During the rollout of the trajectory, the estimator provides the updated belief estimate as an additional feature to the policy. The estimator resets to  $b_0$  at the beginning of each trajectory.

Note that the policy observes the *evolution* of belief along each trajectory. This allows the policy to learn to be optimal with respect to the evolution of belief, and to take actions which affect the

---

**Algorithm 1** Bayes-Adaptive Policy Gradient

---

**Require:** Bayes-Estimator  $ES$ ,  $b_0$ ,  $\mathcal{M}$ , initial policy  $\theta_0$ ,  $n_{\text{itr}}$ ,  $n_{\text{sample}}$

```
1: for  $i = 1, 2, \dots, n_{\text{itr}}$  do
2:   for  $n = 1, 2, \dots, n_{\text{sample}}$  do
3:     Reset  $ES$  with  $b_0$ 
4:     Sample model parameter  $\phi \sim b_0$  to get  $M_\phi$ 
5:     Sample a trajectory  $\tau_n$  from  $M_\phi$  using policy  $\pi(\theta_i)$  and estimator  $ES$ 
6:   Update policy:  $\theta_{i+1} = \text{BatchPolicyOptimization}(\theta_i, \{\tau_1, \dots, \tau_{n_{\text{sample}}}\})$ 
```

---

evolution. For example, given prior belief  $b_0$  with high entropy (i.e., high uncertainty), our algorithm policy would produce conservative or informative actions until  $b$  becomes informative enough along the course of the trajectory. This is a key difference between our algorithm and UP-OSI [8], in which the universal policy has no notion of belief or uncertainty.

## 5 Experiments

Our algorithm falls under the class of model-based algorithms that consider multiple MDPs to train a policy. We compare our algorithm against others from this class, QMDP [4, 5] and EPOpt [7]. QMDP assumes that the latent state  $\phi$  is revealed after one step, taking the action that maximizes  $Q_{\text{MDP}}(s, a, b(\phi)) = \sum_{\phi \in \Phi} b(\phi)Q(s, a, \phi)$ . EPOpt is a minimax-style algorithm which trains the policy to be robust across multiple MDPs, without utilizing the belief over MDPs in learning the policy.

We evaluate the three algorithms on a set of simulated benchmark examples from OpenAI Gym [1] using the MuJoCo physics simulator [13]: Ant, Swimmer, and HalfCheetah. Each of these environments has several model parameters that can be changed to simulate multiple possible MDPs. We uniformly sampled a wide range of parameters (e.g. body link length, mass, joint damping, friction) as described in Table 1. For QMDP and BAPG, which maintain a belief over a discretized set of MDPs, we uniformly sampled a set of  $K = 20$  environments. These  $K$  environments represent the environments that the algorithms have access to in simulation, and can therefore use extensively during training.

We implemented EPOpt with Trust Region Policy Optimization (TRPO) [14] as the underlying batch policy optimization subroutine. We used the implementation of TRPO from `rllab` [15].

QMDP used  $K$  environment-specific policies (trained independently with TRPO on each of the chosen  $K$  MDPs) to estimate  $Q(s, a, \phi)$  with trajectory rollouts. At each time step, each environment-specific policy proposed an action candidate which was evaluated with the QMDP value function. The QMDP agent selected the action candidate with highest  $Q_{\text{MDP}}(s, a, b(\phi))$ .

All the policy networks used in EPOpt and QMDP were Gaussian MLPs, using two hidden layers of 64 units each and tanh activations. Because the BAPG policy network also takes the belief as input, the hidden layers consisted of 128 units and 64 units respectively.

### 5.1 Hypothesis: BAPG policies balance robustness and optimality

We compare the three algorithms on two factors: robustness to model discrepancy and performance. In Figure 3a we plot the average (undiscounted) return computed over 10 rollouts of the learned policies on the unmodified HalfCheetah environment (which was not included as one of the  $K$  environments that were available in simulation, but was the mean of the uniform parameter distributions).

We observe that our algorithm significantly outperforms EPOpt and achieves similar performance to QMDP. We believe that QMDP performs well when the belief quickly becomes certain in a particular MDP, which is close to the true MDP. Then, the QMDP agent will effectively use the policy that was trained on that MDP, which would likely perform well on the true MDP. However, if none of the  $K$  MDPs are close to the true MDP, it is difficult to make performance guarantees for QMDP.

In Figures 3b and 3c, we evaluate the performance of the algorithms on one of the  $K$  MDPs used for training. QMDP clearly outperforms the other two algorithms as it has access to the  $K$  trained

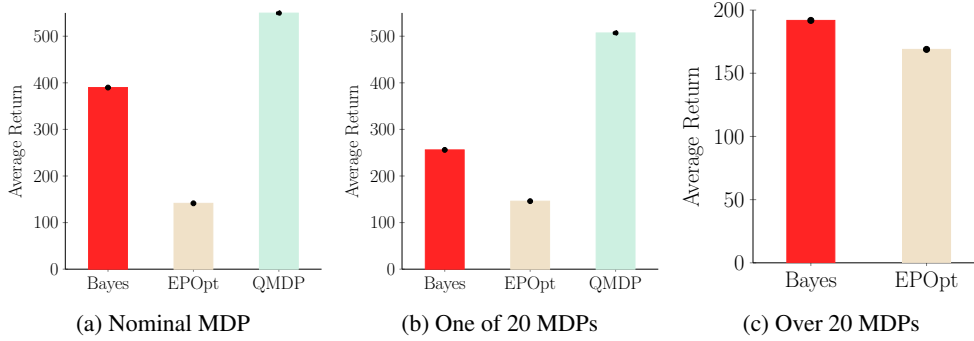


Figure 3: Mean returns on uniformly sampled and prechosen HalfCheetah environments.

policies, and the belief will quickly concentrate around that MDP. We note that our algorithm again outperforms EPOpt on these select MDPs, which is expected since EPOpt does not take advantage of these MDPs during training.

Our experiments demonstrate the robustness of BAPG toward model uncertainty over a wide range of parameter space. We learn one universal policy that adapts to different MDPs by maintaining a belief over the sampled MDPs. In Figure 4, we can see that the sampled MDPs result in drastically different dynamics. Our learned policy network is capable of generating different policies for each of these MDPs to achieve high reward. Videos are available at <https://goo.gl/DN7V9o>.

Parameter	Variation
Damping	$\pm 50\%$
Friction	$\pm 50\%$
Body link size	$\pm 50\%$

Table 1: Model Parameter Variation

## 6 Conclusion and Future Work

We presented a policy-gradient algorithm that, by maintaining a belief over latent physical parameters, achieves robustness to model variations and disturbances without compromising significantly on the optimality. While in this work our belief was over a discrete set of MDPs, there is a natural extension to maintaining a continuous belief distribution either over the MDPs or over the parameter space directly. A possible extension to enable exploration and distinguish policies (and not just beliefs) is to augment the reward to encourage actions that achieve an information gain. An immediately interesting avenue worth exploring is to exploit the mixture of experts framework to consider the action proposals from the MDPs along with the belief.

## References

- [1] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” 2016.
- [2] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, *et al.*, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, no. 7676, p. 354, 2017.
- [3] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen, “Learning hand-eye coordination for robotic grasping with large-scale data collection,” in *International Symposium on Experimental Robotics*, pp. 173–184, Springer, 2016.
- [4] M. L. Littman, A. R. Cassandra, and L. P. Kaelbling, “Learning policies for partially observable environments: Scaling up,” in *Machine Learning Proceedings 1995*, pp. 362–370, Elsevier, 1995.

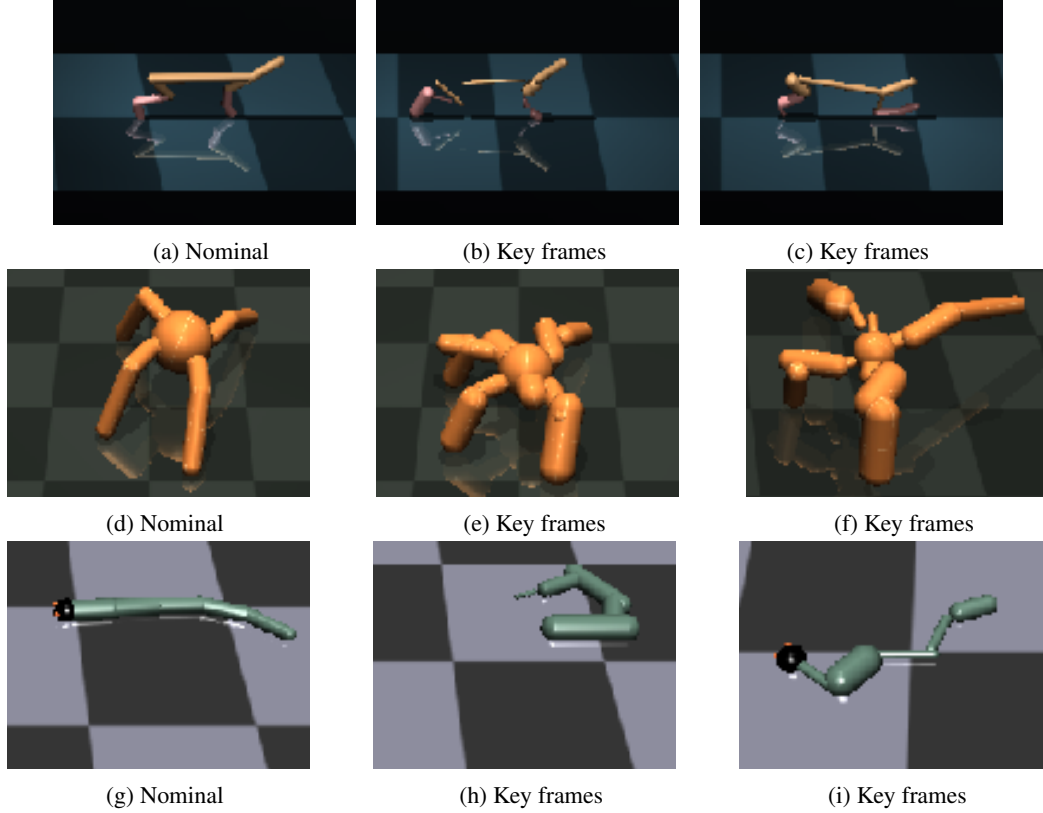


Figure 4: Keyframes from rollouts on various MDPs. Some of the MDPs are geometrically significantly different that they require drastically different policies to achieve optimal returns. See <http://shorturl.at/efxBI> for videos.

- [5] P. Karkus, D. Hsu, and W. S. Lee, “QMDP-Net: Deep Learning for Planning under Partial Observability,” in *Advances in Neural Information Processing Systems*, pp. 4697–4707, 2017.
- [6] M. Chen, E. Frazzoli, D. Hsu, and W. S. Lee, “POMDP-lite for Robust Robot Planning under Uncertainty,” in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*, pp. 5427–5433, IEEE, 2016.
- [7] A. Rajeswaran, S. Ghotra, B. Ravindran, and S. Levine, “EPOpt: Learning Robust Neural Network Policies Using Model Ensembles,” *arXiv preprint arXiv:1610.01283*, 2016.
- [8] W. Yu, J. Tan, C. K. Liu, and G. Turk, “Preparing for the unknown: Learning a universal policy with online system identification,” *arXiv preprint arXiv:1702.02453*, 2017.
- [9] A. Guez, N. Heess, D. Silver, and P. Dayan, “Bayes-adaptive simulation-based search with value function approximation,” in *Advances in Neural Information Processing Systems*, pp. 451–459, 2014.
- [10] M. Ghavamzadeh, S. Mannor, J. Pineau, A. Tamar, *et al.*, “Bayesian Reinforcement Learning: A Survey,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 5-6, pp. 359–483, 2015.
- [11] S. Ross, B. Chaib-draa, and J. Pineau, “Bayes-Adaptive POMDPs,” in *Advances in Neural Information Processing Systems*, pp. 1225–1232, 2008.
- [12] A. Guez, D. Silver, and P. Dayan, “Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search,” in *Advances in Neural Information Processing Systems*, pp. 1025–1033, 2012.

- [13] E. Todorov, T. Erez, and Y. Tassa, “Mujoco: A physics engine for model-based control,” in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 5026–5033, IEEE, 2012.
- [14] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, “Trust region policy optimization,” in *International Conference on Machine Learning*, pp. 1889–1897, 2015.
- [15] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, “Benchmarking Deep Reinforcement Learning for Continuous Control,” in *International Conference on Machine Learning*, pp. 1329–1338, 2016.