

# BA2C: BAYESIAN ADVANTAGE ACTOR CRITIC FOR FEW SAMPLE LEARNING USING FACTOR GRAPH BAYESIAN NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

On-policy reinforcement learning (RL) algorithms, such as Proximal Policy Optimization (PPO), are widely used by researchers and practitioners across various tasks. However, these algorithms are known for their lack of sample efficiency, making them challenging to apply when obtaining training samples is costly, particularly in the absence of an effective simulation environment. While some research exists on Bayesian approaches in the context of RL, which promise a better trade-off between exploration and exploitation, to the best of our knowledge, no prior work has explored the implementation of policy-gradient actor-critic algorithms using expectation-propagation for approximate message passing in Bayesian neural networks (BNNs). In this paper, we propose BA2C, an actor-critic algorithm based on networks represented as a factor graph. Since these networks are trained through approximate message passing rather than gradients, we employ a pseudo-target implementation of the policy gradient theorem. We evaluate our algorithm against three popular RL implementations and observe that required training samples can be reduced up to 50% to reach desired levels on certain environments during the early stages of training. Furthermore, our findings indicate that the uncertainty-based evaluation using expectation propagation actually helps, and that our algorithm performs better within the expectation-propagation approximation compared to IVON, a state-of-the-art variational inference algorithm.

## 1 INTRODUCTION

**Motivation** In recent years, major developments have made the field of reinforcement learning (RL) popular. Although the foundations of RL date back to the 1980s and 1990s, the use of deep neural networks (NNs) has led to major breakthroughs, for example, when Atari games could be solved from pixel input only (Mnih et al., 2013), or when the games of Go (Silver et al., 2017b) and Chess (Silver et al., 2017a) could be solved without human knowledge. Apart from games, RL is heavily used for progress in autonomous driving (Kiran et al., 2021), and it has taken an important role in robot control (Singh et al., 2022).

Robot control typically requires RL agents to handle continuous action spaces, a challenge that only a subset of popular RL algorithms can address. Traditionally, actor critic algorithms like Advantage Actor Critic (A2C) and its successor Proximal Policy Optimization (PPO) are used to solve these kinds of problems. They are on-policy algorithms and trained using the policy gradient theorem. On-policy algorithms are popular because of their reliability and high performance. Especially PPO is one of the most widespread RL algorithms today.

**Goal and Focus of this Paper** However, sample efficiency is a major problem in RL. Even relatively simple problems easily require millions of samples, driving training time and computational effort high. If an accurate and fast simulation environment exists, long training times are inconvenient but manageable. However, if obtaining a large number of samples is expensive or even infeasible, for example, if real-world feedback is needed, the problem of sample efficiency rules out the use of PPO and similar on-policy algorithms. While recent developments have made off-policy algorithms, which use a replay buffer to store and re-use old experiences, more stable, resulting in the development of

054 Soft Actor Critic (SAC) and TD3, the question remains how on-policy algorithms can be made more  
 055 sample-efficient, particularly in early training stages.

057 **Contribution** This paper analyzes whether Bayesian neural networks (BNN) can be used in on-  
 058 policy RL algorithms to improve sample efficiency with better exploration and faster learning. Our  
 059 main contributions are:

- 060 • We propose a variant of the A2C algorithm that is based on neural networks (NNs) repre-  
 061 sented as factor graphs, trained with expectation propagation for message passing (instead  
 062 of normal NNs), which allows for consistent confidence quantification.
- 063 • We adapt the training logic obtained from the policy gradient theorem for the approximate  
 064 message-passing framework that does not use gradients.
- 065 • Our algorithm is applicable to a wide range of continuous control problems.
- 066 • In the implemented algorithm, we compare the training performance side-by-side to standard  
 067 NNs, as they are used in Stable Baselines 3 (Raffin et al., 2021), Ray RLLib (Liang et al.,  
 068 2017), and Dopamine (Castro et al., 2018), three widespread libraries for RL algorithms.
- 069 • For our Bayesian Advantage Actor Critic (BA2C) algorithm, we compare the message-  
 070 passing framework using expectation propagation against IVON, a state-of-the-art variational  
 071 inference (VI) algorithm, cf. Shen et al. (2024).
- 072 • We show that the use of expectation propagation for training BNNs offers great potential to  
 073 improve the sample efficiency of RL algorithms. Our code is available on GitHub<sup>1</sup>  
 074

075 In Section 2, we discuss related work. In Section 3, we propose our approach: BA2C. In Section 4,  
 076 we describe our implementation and design choices. In Section 5, we evaluate our approach against  
 077 other benchmarks. In Section 6, we summarize and discuss the results obtained.

## 079 2 RELATED WORK

081 **RL Methods** RL is thoroughly introduced by Sutton & Barto (2018). RL methods can be catego-  
 082 rized along various dimensions. First, algorithms are split into model-based and model-free; we focus  
 083 on model-free. Second, some methods rely on simple tabular or linear function approximators, while  
 084 others employ deep neural networks for complex tasks. Third, algorithms often differ in whether they  
 085 handle discrete or continuous action spaces.

086 A key distinction is between off-policy and on-policy algorithms. Off-policy methods, such as  
 087 Q-Learning (QL) (Watkins & Dayan, 1992) and Deep Q-Learning (DQN) (Mnih et al., 2015), learn  
 088 from data not necessarily generated by the current policy and typically use a replay buffer. QL was  
 089 originally designed for discrete actions, but Deep Deterministic Policy Gradient (DDPG) (Silver et al.,  
 090 2014) extends QL ideas to continuous actions, with Twin Delayed DDPG (TD3) (Fujimoto et al.,  
 091 2018) and Soft Actor-Critic (SAC) (Haarnoja et al., 2018) improving training stability. By contrast,  
 092 on-policy methods learn strictly from the current policy. They often build on the policy-gradient  
 093 theorem (Sutton et al., 1999). REINFORCE (Williams, 1992) is the earliest such algorithm, followed  
 094 by A2C (Mnih et al., 2016). PPO (Schulman et al., 2017) generalizes A2C and is considered one of  
 095 the most robust current methods (Huang et al., 2022).

096 **Bayesian RL Methods** The area of Bayesian RL has been studied with a variety of methodic  
 097 contributions. As explained in the survey Ghavamzadeh et al. (2015), Bayesian RL offers the  
 098 perspective to solve two classic problems in RL, the trade-off between exploration and exploitation,  
 099 and the option to provide prior knowledge to the agents. The survey gives an overview about several  
 100 approaches on how to incorporate a Bayesian modeling into classical RL algorithms, starting with  
 101 multi-armed bandit settings, and discussing TD learning, SARSA and model-based approaches. It  
 102 also discusses the topic of Partially Observable Markov Decision Processes (POMDPs). However,  
 103 these methods mainly cover tabular cases or linear models. The ability to scale Bayesian RL to more  
 104 complex settings is seen as the major challenge.

105 The survey also discusses an actor-critic algorithm using Bayesian learning, see Ghavamzadeh &  
 106 Engel (2007). In this work, an explicit posterior was derived for simple tabular problems. In a

107 <sup>1</sup><https://github.com/BayesianAdvantageActorCritic/BayesianAdvantageActorCriticSubmission>

subsequent paper Ghavamzadeh et al. (2016), the algorithm was extended by a technique called *Bayesian Quadrature*. More details, proofs, and experimental results were presented. Still, the paper heavily focused on methodic work only with limited scalability.

Achieving higher sample efficiency in exploration was also the motivation for Azizzadenesheli & Anandkumar (2019). The authors utilized Bayesian linear regression for Q-learning. To scale the idea of Bayesian Q-learning to deep neural networks, they used normal neural networks except for the last layer, where they chose Bayesian linear regression instead. They reported to successfully outperform normal deep Q-networks on a variety of Atari games. The idea was taken and adapted, for example, in Ishfaq et al. (2024), Rozanov (2024), and Sasso et al. (2023).

Tasdighi et al. (2024) tackle the common problem of a deficient critic misguiding the actor, resulting in stability problems for actor-critic algorithms. They integrated a probably approximately correct (PAC) Bayesian bound into the SAC algorithm. The algorithm was tested on typical Gym benchmarks, and improvements in sample efficiency were reported, compared to the standard SAC algorithm.

**Factor Graphs & Approximate Message Passing** Our methodological framework is related to factor graphs and approximate message passing techniques using expectation propagation, see, e.g., Yedidia et al. (2005); Minka (2001). These methods — exploiting the sum product algorithm (Aji & McEliece, 2000) together with a Kullback–Leibler-based approximation of the latent variable marginals — allow to efficiently compute belief distributions of latent variables in Bayesian networks and graphical models (Frey, 1998; Wainwright et al., 2008).

Further, we base our work on an approximate message-passing framework for NNs represented as a factor graph (Sommerfeld et al., 2025) that allows to represent the mean and variance of the weights of a standard NN explicitly. Details of this approach are provided in Appendix B.

**Targeted Research Gap** Overall, the use of Bayesian techniques to effectively guide the early exploration of AC-based RL methods is still understudied. To the best of our knowledge, there is no existing approach that (i) is generally applicable, (ii) requires only sparse data (and no sampling techniques), (iii) provides confidence for guided exploration in a natural way, and (iv) works for continuous actions. Our aim is to close that gap.

### 3 BAYESIAN ADVANTAGE ACTOR CRITIC (BA2C)

In this section, we propose our algorithm, an adaptation of advantage actor critic with factor graph BNNs trained via approximate message passing. Our algorithm focuses on problems with a continuous action space. A focused explanation of classical A2C is available in the Appendix A. The core motivation is to improve the exploration behavior by using the uncertainty of the BNN to drive exploration. This uncertainty defines the exploration behavior of our agent during training. The intuition is: If the algorithm is confident with its prediction, this is the case because enough data is collected for this situation, and only little exploration will be done. If there is not much data evidence, the algorithm will be unsure and explore more.

**Setup & Notation** We define our BA2C algorithm for a Markov decision process (MDP) with an  $n$ -dimensional, continuous state space  $\mathcal{S} = \mathbb{R}^n$  and an  $m$ -dimensional, continuous action space  $\mathcal{A} = \mathbb{R}^m$ . On this MDP, the BA2C agent will execute actions that are sampled from a stochastic policy  $\pi$ . The probability density of an action  $a \in \mathcal{A}$  in a state  $s \in \mathcal{S}$  is

$$\pi_s(a) = P(a|s) = \mathcal{N}(a; \mu_s, \text{diag}(\sigma_s^2)) \quad (1)$$

and follows a Gaussian distribution, where  $\mathcal{N}$  is the probability density of a Gaussian distribution. Its parameters  $\mu_s$  and  $\sigma_s$  are determined by the actor. In the multi-dimensional case,  $\mu_s$  and  $\sigma_s$  are vectors, although we only model the individual actions in this action vector to be stochastically independent. Consequently, we write  $\text{diag}(\sigma_s^2)$  for the covariance matrix of the multivariate Gaussian distribution. Due to this independence, the actions can also be calculated separately using BNNs with one-dimensional outputs. For simpler notation, let  $\Pi_S$  be the (Gaussian) random variable describing the action in state  $S$  according to our policy. Like in the standard setup for RL, the policy should maximize the total expected return  $J$  defined as

$$J = \mathbb{E} \left[ \sum_{i=0}^{\infty} r_i \cdot \gamma^i \right] \quad (2)$$

with  $\{S_i\}_{i \in \mathbb{N}}$  as the sequence of visited states, modeled as random variables. For a step  $i$ , the current state  $S_i$  and action  $\Pi_{S_i}$  lead to a (random) reward  $r_i$  and a new state  $S_{i+1}$ . This notation refers to an episode of infinite length, the most general case. Of course, it also works for finite episodes (where all rewards are zero from a certain point).

**A2C with Standard Neural Networks** Appendix A gives a recap on how the classical advantage actor critic works. There, the actor consists of the weights  $w_\mu$  and  $w_\sigma^2$  to calculate two outputs for a state  $s \in \mathcal{S}$ :

$$\mu_{w_\mu} : \mathcal{S} \rightarrow \mathbb{R}^m; \sigma_{w_\sigma} : \mathcal{S} \rightarrow (0, \infty]^m . \quad (3)$$

These outputs are used as parameters (mean and standard deviation) for the Gaussian policy. Although  $\sigma_{w_\sigma}$  is used as the standard deviation of the Gaussian actor, it has no inherent Bayesian semantics. Instead, it is a supporting output that is part of the policy and trained via the policy gradient theorem, as described in equation 15, see Appendix A.

### 3.1 BA2C WITH FG-BNNs

In contrast, our proposed algorithm uses NNs represented as a factor graph (FG-BNNs) as function approximators for both actor and critic. These FG-BNNs also get the state as input, and they are designed to have the same architecture regarding nodes and layers. However, they offer a huge advantage by modelling *distributions* over the function values of actor and critic (where the distribution results from the inherent uncertainty due to a limited training set rather than the uncertainty of the environment only). The factor-graph framework works inherently with Gaussian distributions as approximations, passing natural parameters of normal distributions from layer to layer.

The fact that the framework works with Gaussian distributions is a natural fit to the popular modelling of Gaussian policies for continuous action spaces. The optimal action in a state itself is modeled as an approximate latent variable  $\hat{B}_{\theta_a}(s)$  in the factor graph with parameters  $\theta_a$ . For a state  $s$ , the factor graph calculates a belief on  $\hat{B}_{\theta_a}(s)$ , a Gaussian approximation of its predictive posterior described by a mean vector  $\mu_{\theta_a}(s)$  and the standard deviations of the uncertainty of these means  $\sigma_{\mu_{\theta_a}}(s)$ . Formally, we get a belief distribution for the best action  $a \in \mathcal{A}$  with density,  $s \in \mathcal{S}$ ,

$$f_{\hat{B}_{\theta_a}(s)}(a|s) = \mathcal{N}(a; \mu_{\theta_a}(s), \text{diag}(\sigma_{\mu_{\theta_a}}(s)^2)) . \quad (4)$$

While  $\mu_{\theta_a}$  of these Gaussian output distributions can be viewed as the point-estimates, the standard deviations express the uncertainty resulting from a lack of training data. This standard deviation is small if there is a lot of supporting training data in the state-space area, and it should be large if there is not much training data similar to this sample.

From a high-level perspective, the algorithm appears quite similar to standard A2C (both use a Gaussian policy for a continuous action space with parameters obtained from a parameterized network). But the concept of obtaining the standard deviation to drive the exploration is different.

**BA2C Critic Training** While the critic in case of standard NNs just gives a point estimate, the critic in BA2C is an FG-BNN with parameters  $\theta_c$ . It gives a Gaussian approximation on the state value, representing the uncertainty based on the lack of training data. Let  $\hat{V}_{\theta_c}$  be this random variable. It depends on the state  $s$ , where the probability density is,  $v \in \mathbb{R}$ ,  $s \in \mathcal{S}$ ,

$$f_{\hat{V}_{\theta_c}}(v|s) = \mathcal{N}(v; \mu_{\theta_c}(s), \sigma_{\mu_{\theta_c}}(s)^2) . \quad (5)$$

Here,  $\mu_{\theta_c}(s)$  is the mean parameter and  $\sigma_{\mu_{\theta_c}}(s)$  the standard deviation parameter, that are retrieved from the critic FG-BNN. Note that this variance  $\sigma_{\mu_{\theta_c}}(s)^2$  may not be mixed up with the distribution of actual returns that will be obtained from this state. Instead, it is a measure for the certainty of the estimation of the state value (i.e. a mean). BA2C’s critic is trained via bootstrapping, just like the standard A2C. Hence, for the training of the critic in state  $s$ , we take a sample  $a$  from the approximate, estimated action distribution  $\hat{B}_{\theta_c}$ . This action is executed in the environment, a reward

<sup>2</sup>In most implementations, an actor network shares the same parameters for the mean and variance prediction, except for the last layer.

$r_{s,a}$  is observed, and the new state  $s'$  retrieved. With this information and recursive use of our critic, we calculate a target estimation  $\mathcal{R}(s, a)$  as,  $a \in \mathcal{A}$ ,  $s \in \mathcal{S}$ ,

$$\mathcal{R}(s, a) = r_{s,a} + \gamma V_{\theta_c}(s'). \quad (6)$$

Note that  $\mathcal{R}(s, a)$  is a Gaussian random variable because  $V_{\theta_c}(s)$  is a Gaussian; only linear scaling and addition with real numbers is used. The mean can be easily calculated as

$$\mathbb{E}[\mathcal{R}(s, a)] = \mathbb{E}[r_{s,a} + \gamma V_{\theta_c}(s)] = r_{s,a} + \gamma \mathbb{E}[V_{\theta_c}(s)] = r_{s,a} + \gamma \mu_{\theta_v}(s), \quad (7)$$

and the variance is

$$\mathbb{V}[\mathcal{R}(s, a)] = \mathbb{V}[r_{s,a} + \gamma V_{\theta_c}(s)] = \gamma^2 \mathbb{V}[V_{\theta_c}(s)] = \gamma^2 \sigma_{\mu_{\theta_v}}(s)^2. \quad (8)$$

When training the parameters  $\theta_c$  of the factor graph, these two momentums are set for the regression factor. The continuous calculation with uncertainties takes elegantly account for the fact that the targets are also uncertain, especially during early training.

**Advantages as Random Variables** In a state  $s \in \mathcal{S}$ , let  $a \in \mathcal{A}$  be the action that was taken. Now, we define  $A(a, s) := \mathcal{R}(s, a) - V_{\theta_c}(s)$  as the advantage.<sup>3</sup> It is defined analogously to normal A2C, but with random variables. The interpretation of the reward is if the execution of action  $a$  did improve the situation compared to what is expected in state  $s$ . It is also Gaussian distributed with mean

$$\begin{aligned} \mathbb{E}[A(a, s)] &= \mathbb{E}[\mathcal{R}(s, a) - V_{\theta_c}(s)] = \mathbb{E}[\mathcal{R}(s, a)] - \mathbb{E}[V_{\theta_c}(s)] \\ &= r_{s,a} + \gamma \mu_{\theta_v}(s) - \mu_{\theta_v}(s) = r_{s,a} - \mu_{\theta_v}(s)(1 - \gamma) \end{aligned} \quad (9)$$

due to equation 7. With equation 8 and since we assume independence the variance amounts to

$$\mathbb{V}[A(a, s)] = \mathbb{V}[\mathcal{R}(s, a) - V_{\theta_c}(s)] = \mathbb{V}[\mathcal{R}(s, a)] + \mathbb{V}[V_{\theta_c}(s)] = \gamma^2 \sigma_{\mu_{\theta_v}}(s)^2 + \sigma_{\mu_{\theta_v}}(s)^2. \quad (10)$$

### 3.2 BA2C ACTOR TRAINING VIA PSEUDO TARGETS

An important change from normal A2C to its Bayesian sibling is necessary for the training of the actor. It cannot make direct use of the policy gradient, since the message-passing framework for our FG-BNN does not train via gradients. To tackle this problem, our algorithm makes use of *pseudo targets*. In a state  $s \in \mathcal{S}$ , these pseudo targets  $T_{a,s}$  are obtained as a shift of the currently believed optimal action  $\hat{B}_{\theta_a}(s)$  towards the direction given by the policy gradient theorem, scaled by a learning rate  $\alpha$ . Let  $a$  be the actually taken action which is a sample from the belief distribution, and let  $A(a, s)$  be the advantage.  $\mathcal{N}(a; \mu_{\theta_a}(s), \text{diag}(\sigma_{\mu_{\theta_a}}(s)^2))$  is the probability density of the taken action under this belief, described by the two moments  $\mu_{\theta_a}(s)$  and  $\sigma_{\mu_{\theta_a}}(s)$ . Note that we are only shifting the mean parameter in BA2C. Also, we are taking the derivative just with respect to the mean  $\mu_{\theta_a}$ , not the parameters  $\theta_a$ . Then, the pseudo targets  $T_{a,s}$  are defined as,  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ ,

$$\begin{aligned} T_{a,s} &= \hat{B}_{\theta_a}(s) + \alpha \cdot A(a, s) \cdot \nabla_{\mu_{\theta_a}} \log [\mathcal{N}(a; \mu_{\theta_a}(s), \text{diag}(\sigma_{\mu_{\theta_a}}(s)^2))] \\ &= \hat{B}_{\theta_a}(s) + \alpha \cdot A(a, s) \cdot \nabla_{\mu_{\theta_a}} \log \left[ \frac{1}{\sqrt{2\pi\sigma_{\mu_{\theta_a}}(s)}} \cdot \exp \left[ -\frac{(a - \mu_{\theta_a}(s))^2}{2\sigma_{\mu_{\theta_a}}(s)^2} \right] \right] \\ &= \hat{B}_{\theta_a}(s) + \alpha \cdot A(a, s) \cdot (a - \mu_{\theta_a}(s)) / (2\sigma_{\mu_{\theta_a}}(s)^2). \end{aligned} \quad (11)$$

Fortunately, the derivative of the log probability density of a normal distribution with respect to the mean parameter is quite simple. Because  $\hat{B}_{\theta_a}(s)$  and  $A(a, s)$  are Gaussian random variables, the other variables are constants, and Gaussian distributions are closed under linear scaling and addition,  $T_{a,s}$  is also a Gaussian random variable. We can calculate the mean parameter vector  $\mathbb{E}[T_{a,s}]$  in a straightforward way with the results from equation 11 and equation 9:

$$\begin{aligned} \mathbb{E}[T_{a,s}] &= \mathbb{E} \left[ \hat{B}_{\theta_a}(s) + \alpha \cdot A(a, s) \cdot \frac{a - \mu_{\theta_a}(s)}{2\sigma_{\mu_{\theta_a}}(s)^2} \right] = \mathbb{E}[\hat{B}_{\theta_a}(s)] + \alpha \cdot \mathbb{E}[A(a, s)] \cdot \frac{a - \mu_{\theta_a}(s)}{2\sigma_{\mu_{\theta_a}}(s)^2} \\ &= \mu_{\theta_a}(s) + \alpha \cdot (r_{s,a} - \mu_{\theta_v}(s)(1 - \gamma)) \cdot (a - \mu_{\theta_a}(s)) / (2\sigma_{\mu_{\theta_a}}(s)^2). \end{aligned} \quad (12)$$

<sup>3</sup> $V_{\theta_c}(s)$  and  $\mathcal{R}(s, a)$  as defined in the previous section

With the reward  $r_{s,a}$ , this can be explicitly calculated in the training loop. Further, with zero covariance due to the assumed independence between the belief distribution and the advantage distribution, we get a very similar calculation for the variance:

$$\begin{aligned} \mathbb{V}[T_{a,s}] &= \mathbb{V}\left[\hat{B}_{\theta_a}(s) + \alpha \cdot A(a,s) \cdot \frac{a - \mu_{\theta_a}(s)}{2\sigma_{\mu_{\theta_a}}(s)^2}\right] = \mathbb{V}\left[\hat{B}_{\theta_a}(s)\right] + \alpha^2 \cdot \mathbb{V}\left[A(a,s)\right] \cdot \left(\frac{a - \mu_{\theta_a}(s)}{2\sigma_{\mu_{\theta_a}}(s)^2}\right)^2 \\ &= \sigma_{\mu_{\theta_a}}(s)^2 + \alpha^2 \cdot (\gamma^2 \sigma_{\mu_{\theta_v}}(s)^2 + \sigma_{\mu_{\theta_v}}(s)^2) \cdot ((a - \mu_{\theta_a}(s))/(2\sigma_{\mu_{\theta_a}}(s)^2))^2, \end{aligned} \quad (13)$$

using equation 11 and equation 10. These two parameters define the training target for the training of the actor.

## 4 IMPLEMENTATION

In this section, we give an overview of the most important aspects of the implementation. Detailed explanations of the design choices are provided in Appendix C.1. The BA2C algorithm is implemented in Julia. In our BA2C implementation, the FG-BNNs are accessed via a light API. This enables us to exchange FG-BNNs for other types of BNNs or standard NNs while keeping the rest of the BA2C algorithm, including the pseudo-target calculation, unchanged. Further, this allows us to distinguish between the effects of the FG-BNNs and those of the design decisions for BA2C. We are making use of this in Sections 5.1, 5.2, and 5.3. In Appendix C.2, we explain a few design choices for greater stability, such as independence for one training batch, clipping of the pseudo targets, and decreasing the learning rate. Review Appendix C.3 for notable differences in comparison to standard NNs and Appendix C.4 for an explanation of new hyperparameters that come with the selection of FG-BNNs. There, a comprehensive overview of the selected hyperparameters is given in Table 1.

## 5 EXPERIMENTS & EVALUATION

We use different experiments to evaluate the strengths and weaknesses of our Bayesian approach. Within our implementation of the actor critic agent, we are going to evaluate the performance against normal NNs and IVON to justify the choice of the message-passing framework. Afterwards, we will compare the implementation against three of the most widespread implementations of PPO.

**Evaluation Methods** The focus of our evaluation is to assess if the Bayesian approach has an advantage in the early phasis of training. To discuss that, we are using learning curves, a standard way of evaluating RL algorithms. For a given time step  $t$ , the learning curve  $L(t)$  gives an estimate of the average per-step reward of the model after  $t$  timesteps were trained. It is calculated as an exponentially moving average with the previous rewards, calculated recursively as  $L(t) = (1 - \theta)L(t - 1) + \theta r_t$ , where  $r_t$  is the reward after the  $t$ -th step has been executed.<sup>4</sup>  $\theta$  is a coefficient between 0 and 1, indicating the amount of smoothing on the pure rewards. We selected  $\theta = 0.0003$ .

Of course, for a recursive definition, we need to give a base case  $L(0)$ . For a specific problem, we introduce a default base case, the reward that an untrained agent usually yields. We use the same base case for all algorithms we compare on a specific environment. We decided to not just take the first reward as the base case, because the rewards have a high variance and our low  $\theta$  puts heavy emphasis on the base case, especially in the early phasis of the training. For each configuration, we train 5 agents independently. In our plots, for a given time step, we show the interval between the minimum and maximum with a transparent color, and highlight the median with a line of the same color.

**Gym & Pendulum Environment** The implementation was evaluated on several widespread standard-environments given by Gymnasium, the successor of OpenAI Gym (Brockman et al., 2016). Because the implementation is designed for continuous action spaces, we focus on standard and MuJoCo environments. First, we take a deep-dive with the `Pendulum-v1` environment. `Pendulum-v1` is an easy, but non-trivial environment with a three-dimensional observation space

<sup>4</sup>We do not use the cumulated episode reward because we are using a large number of environments, while only using a limited number of samples from each environment. Therefore, we do not have enough data to reliably estimate the episode lengths, especially in the early phasis, when many environments terminate early.

and a one-dimensional action space. The actions apply a torque on a pendulum to bring and keep it in an upright position. Negative reward is given depending on the angle that is missing towards the upright position and for higher absolute rotation speed. Ideally, the pendulum stays perfectly upright at zero speed during the whole episode. This would result in a perfect reward of zero. However, due to random initialization, the agent first needs to swing the pendulum into this position. This process is a nice benchmark for the trade-off between fast and long-term reward.

## 5.1 DIRECT COMPARISON BETWEEN FACTOR GRAPH BNNs AND PyTORCH NEURAL NETWORKS

First, we want to benchmark our algorithm against itself with standard NNs. As explained in Appendix C.1, BA2C can be operated with standard NNs. The whole logic is identical to the case with BNNs, including the calculation of the pseudo targets. We implemented an actor and critic via PyTorch (Paszke et al., 2019), which we use in our algorithm instead of the factor graph framework. The chosen architecture is the same like in Stable Baselines 3 (SB3), with two hidden layers of 64 neurons each, and tanh-activation function in between the layers.<sup>5</sup> As standard NNs do not propagate a standard deviation, the Python implementation usually returns  $\beta_{actor}$  and  $\beta_{critic}$  as constants.

Standard NNs and FG-BNNs can be set up to four combinations: (i) standard for both actor and critic, (ii) standard actor and Bayesian critic, (iii) vice versa, and (iv) both Bayesian. These four combinations were evaluated in Figure 1(a). Clearly, the agent improves its performance the fastest, if the factor graph NNs are used for both actor and critic (blue). If the actor gets a normal NN instead, the learning gets significantly slower, taking around the double amount of samples to reach a level of -2 (yellow). However, the training is still very stable, maxing out at practically the same level as the option with both NNs. If standard NNs are used for both actor and critic, the training gets significantly slower and less stable (green), eventually reaching the same performance on a longer training horizon (Figure 9 in Appendix E). Interestingly, the variant with a BNN as actor and a normal NN as critic does not improve to an acceptable level (green).

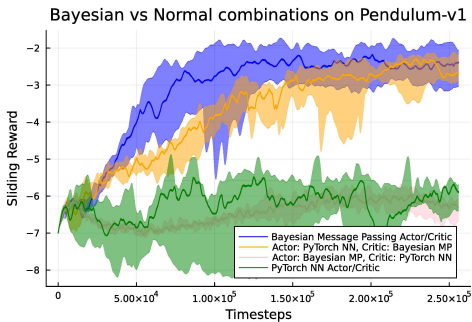
The comparison demonstrates that the BNN is more important for the critic than the actor, a phenomena worth investigating. We can explain the observation by plotting the mean squared error (MSE) of the critic network during training. As shown in Appendix E, the MSE is orders of magnitudes higher in both configurations with a normal NN as critic. The performance of the training is bottlenecked by the speed of the critic. That underlines the fact that the factor graph NNs learn faster, resulting in a lower error given the current policy. If the actor is also a normal NN, actor and critic learn at the same speed, resulting in slower, but still stable learning. If the actor is a BNN, the actor virtually “runs away” from the critic. As a result, the critic’s value estimation is not up-to-date with the actor’s policy, resulting in instability. The opposite situation, where the critic is Bayesian and the actor is normal, is stable, but slower. Notably, variants with a normal actor have a better asymptotic performance. They outperform the Bayesian variant at a certain point, but take longer to achieve this performance.

## 5.2 DOES THE BNN EXPLORATION MECHANISM HELP?

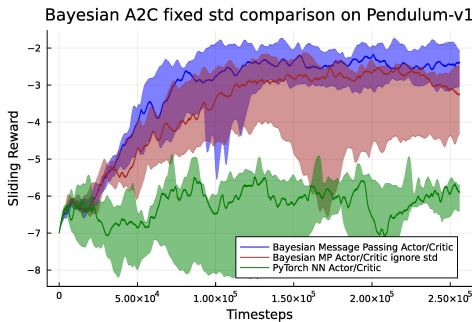
The previous experiment shows that the use of FG-BNNs indeed improves the training speed on the Pendulum environment. It is an interesting question, whether this improvement in training speed arises from the uncertainty-driven exploration of the BNNs compared to the fixed standard deviation which was used for the normal NNs, or if other properties are responsible. To test that, we set up a new experiment in which we benchmark three configurations: First, we use our BA2C agent in standard configuration (both actor and critic with FG-BNNs and uncertainty-driven exploration). Second, we use a BA2C, but instead of using the standard deviation obtained from the networks, we fix the standard deviations to  $\beta_{actor}$ , and  $\beta_{critic}$  for actor and critic. Third, we use the configuration with standard NNs and fixed standard deviation as in the previous experiment. As we can see in Figure 1(b), learning speed is clearly better if the uncertainties are used for exploration. Nevertheless, BNNs still outperform normal NNs when no uncertainty-driven exploration happens.

<sup>5</sup>To allow exact comparability to the Bayesian implementation, the last layer of the actor does not squeeze the output into  $[-1, 1]$  via tanh, but uses clipping as well. Check Appendix C.2.

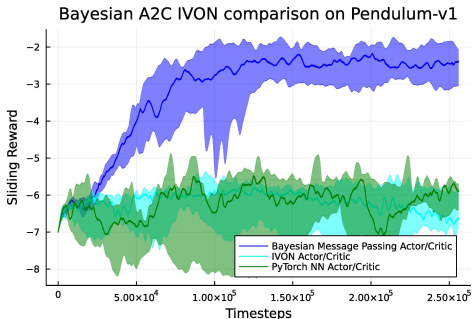
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431



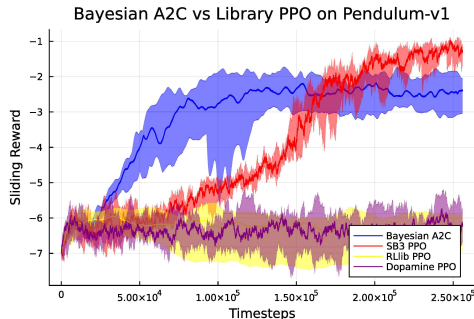
(a) This plot shows the learning curves of the RL agents on the Pendulum-v1 environment over 256,000 training steps. For both the actor and the critic, a standard NN and an FG-BNN can be used, resulting in four combinations; Section 5.1.



(b) Comparison of learning curves for Bayesian Actor Critic BA2C (blue), Bayesian Actor Critic BA2C with *fixed* standard deviation (red) and standard NNs (green); Section 5.2.



(c) The performance of our algorithm with IVON (turquoise) for actor and critic against our standard configuration (blue) and standard NNs (green); Section 5.3.



(d) Comparison of our Bayesian actor critic BA2C with FG-BNNs (blue) against the PPO implementation of SB3 (red), Ray RLLib (yellow), and Google’s Dopamine (purple); Section 5.4.

Figure 1: Performance Comparison of variants of BA2C, cf. (a)-(b), against IVON (c) and PPO (d) for Pendulum; for seven further Gymnasium environments with similar results, see Appendix F.

### 5.3 COMPARISON TO IVON

The choice of the factor graph framework to implement the agent is not straightforward, as there are several approaches for integrating uncertainty estimation into NNs. One widespread family of algorithms is variational inference (Zhang et al., 2018). Here, state of the art is IVON, showing significant advances even for large-scale NNs (Shen et al., 2024). Hence, IVON is a strong candidate for Bayesian RL as well, and it is interesting to test how the algorithm performs when trained with IVON for actor and critic. Here, we followed the hyperparameter guide of the IVON paper. With IVON for actor and critic, cf. Figure 1(c), the algorithm does not improve for Pendulum. In Appendix D, we did an ablation study over hyperparameters of IVON and found no configuration suitable. However, for other environments, we found that IVON indeed performed well (Figure 13).

### 5.4 LIBRARY COMPARISON

For a holistic comparison, the performance of our agents must be compared to state-of-the-art library implementations of policy-gradient actor-critic algorithms. For our comparison, we selected PPO implementations of three of the most popular RL libraries. SB3, cf. Raffin et al. (2021), is a widespread library and known for its high performance. The implementation is based on PyTorch. Because it is a direct successor of the original PPO implementation in the *Baselines*-library (Dhariwal et al., 2017), which has been unsupported since 2020, it is a good candidate for comparison. Ray RLLib (Liang et al., 2017) claims to be a “Industry-Grade, Scalable Reinforcement Learning” library, supporting distributed learning on large clusters. Ray RLLib supports both TensorFlow and PyTorch, but prefers PyTorch in their recent versions. Hence, we are using the PyTorch configuration.

432 Dopamine (Castro et al., 2018) is an RL framework open-sourced by Google and positions itself  
 433 as “a research framework for fast prototyping of RL algorithms”. It is based on TensorFlow. These  
 434 implementations were three of the RL libraries with the most stars on GitHub (9.5k for SB3, 34.7k  
 435 for the Ray framework, and 10.6k for Dopamine). So, they are a diverse and powerful group of  
 436 competitors for our BA2C implementation. We ran all library algorithms with their given standard-  
 437 configuration, but set  $\gamma$  to 0.9 and the batch size to 256 to ensure same conditions for the library and  
 438 our BA2C. On Pendulum, we verified that the selection of 0.9 instead of 0.99 for  $\gamma$  improved the  
 439 training speed of the standard PPO, and did not lower the final performance on the short horizon  
 440 of early training. Figure 1(d) shows that especially SB3’s PPO is a stable and reliable algorithm to  
 441 solve Pendulum. It reaches the level of -3 after around 160,000 steps, which is around double the  
 442 amount compared to the median of BA2C. agent, demonstrating a strong advantage of the Bayesian  
 443 agent during early training. When looking at the level of -4, the advantage of BA2C is even stronger,  
 444 where SB3 needs up to three times as many samples. However, in the later phasis of training, SB3  
 445 outperforms our Bayesian Agent in terms of peak-performance. Ray RLlib and Dopamine mostly fail  
 446 to improve compared to the initial level.

447 Of course, our algorithm was tested on more environments than just Pendulum. The results are  
 448 printed in the Appendix F. We will summarize the findings here. While SB3 often performs best in  
 449 terms of final performance, the Bayesian algorithm and Ray RLlib are superior when it comes to  
 450 early learning. The race between BA2C and Ray RLlib is close, with the Bayesian being faster on  
 451 some environments, and Ray RLlib being faster on others. However, Ray RLlib is unable to solve  
 452 Pendulum, and there is no environment where Ray RLlib clearly outperforms the Bayesian agent.  
 453 Dopamine falls behind the other implementations.

## 454 6 LIMITATIONS, FUTURE WORK & CONCLUSION

456 **Limitations** The actor critic algorithm with BNNs demonstrates stable and fast learning. However,  
 457 the training speed was significantly reduced, while stability is also affected, and hyperparameter  
 458 sensitivity is increased. Careful tuning, for example, with the GAE and reward and advantage  
 459 normalization could help mitigating this problem. Nevertheless, a lower  $\gamma$  does not hurt the value  
 460 proposition of this algorithm, since it is designed for cases with sparse environment access, situations  
 461 that often do not allow for a long lookahead.

462 The training with BNNs is more sample efficient, but is very time-consuming due to the high  
 463 computation demand of the BNN framework. The training part of the algorithm takes around 50 times  
 464 longer than with a PyTorch implementation. The longer computation time has two main reasons.  
 465 First, the Bayesian frameworks, including the one we use, are results of recent research projects. In  
 466 comparison to the highly optimized and professionally used PyTorch, they focus on a conceptual  
 467 implementation instead of compute efficiency. Second, the probabilistic computations are just more  
 468 complex mathematically, requiring more operations. Hence, a comparison with PyTorch in terms of  
 469 resource efficiency is not expedient. Nevertheless, we think the optimization of Bayesian frameworks  
 470 for better efficiency pose a worthy research opportunity.

471 **Future Work** Although the Bayesian implementation provides significant improvements to the  
 472 early training performance of on-policy algorithms, the sample efficiency of off-policy algorithms like  
 473 SAC can be superior. In future research, one could utilize BNNs in the SAC algorithms to combine  
 474 the benefits with regard to sample efficiency.

476 **Conclusion** We proposed a variant of the A2C that works with BNNs which are not trained via  
 477 gradients. Instead, we developed a pseudo-target formulation to train the actor. The use of BNNs  
 478 improves the sample efficiency over standard NNs within the same algorithm significantly. When  
 479 comparing the approach to highly mature RL implementations like SB3, the algorithm can prevail  
 480 when it comes to the number of necessary samples to achieve competitive performance, and even  
 481 outperforms the algorithm in the early phases.

482 The high training speed in the number of samples offers great opportunity for practitioners who need  
 483 to work with sparse samples. For example, in robot control tasks, there is sometimes no accurate  
 484 simulation environment available, and training is done with the real robot. In this case, the time  
 485 required to optimize the parameters is minimal compared to the time for obtaining the samples. Here,  
 our Bayesian approach can significantly improve training times.

## REFERENCES

- 486  
487  
488 Srinivas M Aji and Robert J McEliece. The generalized distributive law. *IEEE Transactions on*  
489 *Information Theory*, 46(2):325–343, 2000.
- 490 Kamyar Azizzadenesheli and Animashree Anandkumar. Efficient exploration through bayesian deep  
491 q-networks, 2019. URL <https://arxiv.org/abs/1802.04412>.
- 492  
493 Lukas Biewald. Experiment tracking with weights and biases, 2020. URL [https://www.wandb.](https://www.wandb.com/)  
494 [com/](https://www.wandb.com/). Software available from wandb.com.
- 495  
496 Greg Brockman et al. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 497 Pablo Samuel Castro, Subhodeep Moitra, Carles Gelada, Saurabh Kumar, and Marc G. Bellemare.  
498 Dopamine: A research framework for deep reinforcement learning. *CoRR*, abs/1812.06110, 2018.  
499 URL <http://arxiv.org/abs/1812.06110>.
- 500  
501 Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford,  
502 John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. [https:](https://github.com/openai/baselines)  
503 [//github.com/openai/baselines](https://github.com/openai/baselines), 2017.
- 504  
505 Brendan J Frey. *Graphical models for machine learning and digital communication*. MIT press,  
506 1998.
- 507  
508 Scott Fujimoto et al. Addressing function approximation error in actor-critic methods. *CoRR*,  
509 abs/1802.09477, 2018.
- 510  
511 Mohammad Ghavamzadeh and Yaakov Engel. Bayesian actor-critic algorithms. In *Proceedings of*  
512 *the 24th International Conference on Machine Learning, ICML '07*, pp. 297–304, New York, NY,  
513 USA, 2007. Association for Computing Machinery. ISBN 9781595937933. doi: 10.1145/1273496.  
1273534. URL <https://doi.org/10.1145/1273496.1273534>.
- 514  
515 Mohammad Ghavamzadeh, Yaakov Engel, and Michal Valko. Bayesian policy gradient and actor-  
516 critic algorithms. *Journal of Machine Learning Research*, 17(1):2319–2371, January 2016. ISSN  
1532-4435.
- 517  
518 Mohammed Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Convex optimization:  
519 Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(5–6):359–483,  
520 2015. ISSN 1935-8245. doi: 10.1561/22000000049. URL [http://dx.doi.org/10.1561/](http://dx.doi.org/10.1561/22000000049)  
521 [22000000049](http://dx.doi.org/10.1561/22000000049).
- 522  
523 Tuomas Haarnoja et al. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning  
524 with a stochastic actor. In *ICML 2018, 10-15, 2018*, volume 80 of *Proceedings of Machine*  
*Learning Research*, pp. 1856–1865. PMLR, 2018.
- 525  
526 Shengyi Huang, Anssi Kanervisto, Antonin Raffin, Weixun Wang, Santiago Ontañón, and Rousslan  
527 Fernand Julien Dossa. A2C is a special case of PPO. <https://arxiv.org/abs/2205.09123>, 2022.
- 528  
529 Haque Ishfaq, Qingfeng Lan, Pan Xu, A. Rupam Mahmood, Doina Precup, Anima Anandkumar, and  
530 Kamyar Azizzadenesheli. Provable and practical: Efficient exploration in reinforcement learning  
via langevin monte carlo, 2024. URL <https://arxiv.org/abs/2305.18246>.
- 531  
532 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil  
533 Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey,  
534 2021. URL <https://arxiv.org/abs/2002.00444>.
- 535  
536 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Joseph Gonzalez, Ken  
537 Goldberg, and Ion Stoica. Ray rllib: A composable and scalable reinforcement learning library.  
*CoRR*, abs/1712.09381, 2017. URL <http://arxiv.org/abs/1712.09381>.
- 538  
539 Thomas Peter Minka. *A family of algorithms for approximate Bayesian inference*. PhD thesis,  
Massachusetts Institute of Technology, 2001.

- 540 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
541 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013. URL  
542 <https://arxiv.org/abs/1312.5602>.  
543
- 544 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim  
545 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement  
546 learning, 2016. URL <https://arxiv.org/abs/1602.01783>.
- 547 Volodymyr Mnih et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):  
548 529–533, 2015. ISSN 14764687.
- 549
- 550 Adam Paszke et al. Pytorch: An imperative style, high-performance deep learning library, 2019.
- 551
- 552 Antonin Raffin et al. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*  
553 *Machine Learning Research*, 22(268):1–8, 2021.
- 554 Nikolai Rozanov. Efficient exploration in deep reinforcement learning: A novel bayesian actor-critic  
555 algorithm, 2024. URL <https://arxiv.org/abs/2408.10055>.
- 556
- 557 Remo Sasso, Michelangelo Conserva, and Paulo Rauber. Posterior sampling for deep reinforcement  
558 learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan  
559 Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine*  
560 *Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 30042–30061. PMLR,  
561 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/sasso23a.html>.
- 562 John Schulman et al. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*,  
563 2017.
- 564
- 565 Yuesong Shen, Nico Daheim, Bai Cong, Peter Nickl, Gian Maria Marconi, Clement Bazan, Rio  
566 Yokota, Iryna Gurevych, Daniel Cremers, Mohammad Emtiyaz Khan, and Thomas Möllenhoff.  
567 Variational learning is effective for large deep networks, 2024. URL <https://arxiv.org/abs/2402.17641>.  
568
- 569 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,  
570 Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Si-  
571 monyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement  
572 learning algorithm, 2017a. URL <https://arxiv.org/abs/1712.01815>.
- 573
- 574 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,  
575 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go  
576 without human knowledge. *Nature*, 550(7676):354–359, 2017b. doi: 10.1038/nature24270. URL  
577 <https://doi.org/10.1038/nature24270>.
- 578 David Silver et al. Deterministic policy gradient algorithms. In *ICML’14, Vol. I*, pp. 387–395, 2014.
- 579
- 580 Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. Reinforcement learning in robotic ap-  
581 plications: a comprehensive survey. *Artif. Intell. Rev.*, 55(2):945–990, February 2022. ISSN  
582 0269-2821. doi: 10.1007/s10462-021-09997-9. URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10462-021-09997-9)  
583 [s10462-021-09997-9](https://doi.org/10.1007/s10462-021-09997-9).
- 584 Leslie N. Smith. A disciplined approach to neural network hyper-parameters: Part 1 – learning rate,  
585 batch size, momentum, and weight decay, 2018. URL [https://arxiv.org/abs/1803.](https://arxiv.org/abs/1803.09820)  
586 [09820](https://arxiv.org/abs/1803.09820).
- 587
- 588 Romeo Sommerfeld, Christian Helms, and Ralf Herbrich. Approximate message passing for bayesian  
589 neural networks, 2025. URL <https://arxiv.org/abs/2501.15573>.
- 590
- 591 David Stern, Ralf Herbrich, and Thore Graepel. Matchbox: Large scale online bayesian recommen-  
592 dations. pp. 111–120, 04 2009. doi: 10.1145/1526709.1526725.
- 593
- 593 Richard S. Sutton and Andrew G. Barto. *Reinforcement learning - An introduction*. Adaptive  
computation and machine learning. MIT Press, second edition, 2018. ISBN 978-0-262-19398-6.

594 Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods  
595 for reinforcement learning with function approximation. In *Proceedings of the 13th International*  
596 *Conference on Neural Information Processing Systems*, NIPS'99, pp. 1057–1063, Cambridge, MA,  
597 USA, 1999. MIT Press.

598 Bahareh Tasdighi, Abdullah Akgül, Manuel Hausmann, Kenny Kazimirzak Brink, and Melih  
599 Kandemir. Pac-bayesian soft actor-critic learning, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2301.12776)  
600 [2301.12776](https://arxiv.org/abs/2301.12776).

601 Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational  
602 inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.

603 Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992.

604 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
605 learning. *Mach. Learn.*, 8(3–4):229–256, May 1992. ISSN 0885-6125. doi: 10.1007/BF00992696.  
606 URL <https://doi.org/10.1007/BF00992696>.

607 J.S. Yedidia, W.T. Freeman, and Y. Weiss. Constructing free-energy approximations and generalized  
608 belief propagation algorithms. *IEEE Transactions on Information Theory*, 51(7):2282–2312, 2005.  
609 doi: 10.1109/TIT.2005.850085.

610 Cheng Zhang, Judith Butepage, Hedvig Kjellstrom, and Stephan Mandt. Advances in variational  
611 inference, 2018. URL <https://arxiv.org/abs/1711.05597>.

612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## 648 A BACKGROUND: ADVANTAGE ACTOR CRITIC

649  
650 A widespread RL algorithm is Advantage Actor Critic (A2C). This is an on-policy algorithm, which  
651 means that only samples obtained from the current policy can be used to train the agent. A2C uses  
652 two kinds of neural networks (NNs), an actor and a critic. Let us discuss the actor first, and motivate  
653 the critic later.

### 654 A.1 ACTOR & POLICY GRADIENT THEOREM

655 For a given Markov decision process (MDP), let  $\mathcal{S} = \mathbb{R}^n$  be an  $n$ -dimensional state space, and  
656  $\mathcal{A} = \mathbb{R}^m$  an  $m$ -dimensional action space. While A2C works both with discrete and continuous action  
657 spaces, we will entirely focus on continuous problems in this work. Actor Critic algorithms work  
658 with stochastic policies, so the policy  $\pi$  of the RL agent gives a probability density function over the  
659 action space, given a current state:

$$660 \pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, \infty). \quad (14)$$

661 The actions applied in the environment in a given state are sampled from this distribution. In the  
662 vast majority of implementations for the continuous case, a Gaussian distribution is used for the  
663 policy function. Hence, the actor is an NN with weights  $w$  that takes a state and outputs a mean and a  
664 standard deviation for each action dimension. In the multidimensional case, the actions are usually  
665 predicted independently (all covariances zero). We will write the means as  $\mu_w(s)$  and the standard  
666 deviations as  $\sigma_w(s)$ .

667 The idea behind the training is intuitively easy to understand. During the training, actions are sampled,  
668 executed, and the reward is measured. If, in a specific state  $s$  – a specific action  $a$  turns out to be  
669 *good*, the probability is increased. If the action turns out to be *bad*, the probability is decreased. To  
670 increase or decrease the probability of  $a$ , the gradient of the probability of  $a$  with respect to  $w$  will be  
671 calculated, and gradient ascent is used to move the weights in the direction of the gradient.

672 The theoretical foundation for this is given by the policy gradient theorem from the literature. Let  
673  $J(w)$  denote the expected (discounted) return when following policy  $\pi_w$  (for example, averaged over  
674 all relevant starting states). The policy gradient theorem states that

$$675 \nabla_w J(w) = \mathbb{E}_{s \sim d^\pi, a \sim \pi_w(\cdot | s)} \left[ A(a, s) \nabla_w \log(\pi_w(a | s)) \right], \quad (15)$$

676 where  $d^\pi$  is the distribution of states under policy  $\pi_w$ , and  $A(a, s)$  is the advantage of action  $a$  in  
677 state  $s$ . In practice, we approximate this expectation by a single transition (or a batch of transitions),  
678 and perform a gradient ascent update of the form

$$679 w \leftarrow w + \alpha A(a, s) \nabla_w \log(\mathcal{N}(a; \mu_w(s), \text{diag}(\sigma_w(s)^2))). \quad (16)$$

680 Here,  $\mathcal{N}(a; \mu_w(s), \sigma_w(s)^2)$  is the Gaussian probability density for action  $a$  under the current policy  
681 parameters  $w$ , and  $\alpha$  is the learning rate. This formula makes it clear that the probability of actions  
682 with positive advantage is increased, whereas it is decreased for actions with negative advantage.

### 683 A.2 REWARD CALCULATION AND CRITIC

684 It is closely related to the reward  $r$  the step gave, but it is more than that. The advantage measures  
685 the difference between the value (expected return) of the next state  $V(s')$  and the current state  $s$   
686 (discounted):

$$687 A(a, s) = (r + \gamma V(s')) - V(s). \quad (17)$$

688 How do we know the state value? Here comes the critic into play. The critic is a neural network that  
689 calculates the function

$$690 V : \mathcal{S} \rightarrow \mathbb{R} \quad (18)$$

691 that maps a state to an estimate of its state value, the expected discounted sum of rewards from this  
692 step onward. It gets trained via bootstrapping, which means that predictions from the network itself  
693 for the subsequent step will be used for the update. Mathematically, the target for the updated state  
694 value  $V^*(s)$  is

$$695 V^*(s) = r + \gamma V(s'). \quad (19)$$

696 With these formulas, we have everything needed to understand the vanilla version of Advantage Actor  
697 Critic (A2C).

### 702 A.3 EXPLORATION AND TRAINING

703  
704 Finally, we need to discuss how exploration happens in A2C. Like every RL algorithm, it needs to  
705 tackle the trade-off between exploration and exploitation. How much should the algorithm focus to  
706 learn more about what it has already discovered, and how much should it try something else? Actor  
707 Critic does this in a quite elegant way: The stochastic policy assigns high probability to its preferred  
708 actions (the mean of the Gaussian distribution), but it always gives some probability to completely  
709 other actions. How much exploration happens is determined by the standard deviation parameter,  
710 which is an output of the NN. It gets trained via the gradient, just like the mean. During training, four  
711 cases can happen that are important to understand on an intuitive level:

- 712 1. The taken action is close to the mean (distance less than one standard deviation), and the  
713 advantage is positive: In this case, the training direction is to reduce the standard deviation,  
714 since this increases the probability density within those standard deviations.
- 715 2. The taken action is outside the range between the two inflection points and the advantage is  
716 positive: The standard deviation will be increased since this raises the probability density  
717 for the taken action.
- 718 3. If the action is inside the two inflection points and the advantage is negative, the standard  
719 deviation will be increased.
- 720 4. In the remaining case, the standard deviation will be decreased.

721  
722 This change of the standard deviation is by no means arbitrary. It fully aligns with the policy gradient  
723 theorem, but it can be read as an auxiliary concept next to the desired training of the mean. How  
724 much exploration happens in a specific step does not have a real semantics, and regularly causes  
725 problems in practical applications. To tackle these problems, regularization is needed, for example an  
726 additional entropy loss in SAC.

## 727 B BACKGROUND: APPROXIMATE MESSAGE PASSING FOR BNNs

728  
729 The core objective in BNNs is to determine the predictive posterior distribution  $p(y | x, \mathcal{D})$ , repre-  
730 senting the probability of an output  $y$  given an input  $x$  and a training dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$   
731 of independent and identically distributed samples  $(x_i, y_i)$ . This distribution can be theoretically  
732 obtained by integrating over the network’s weights  $\theta$ , which are treated as random variables with  
733 prior beliefs  $p(\theta)$ :

$$734 \quad p(y | x, \mathcal{D}) = \int p(y | x, \theta) p(\theta | \mathcal{D}) d\theta. \quad (20)$$

735  
736 The posterior distribution of the weights,  $p(\theta | \mathcal{D})$ , is proportional to the product of the prior and the  
737 likelihood of  $\mathcal{D}$ , i.e.,

$$738 \quad p(\theta | \mathcal{D}) \propto p(\theta) \prod_{i=1}^n p(y_i | f_{\theta}(x_i)), \quad (21)$$

739 where  $f_{\theta}$  is the function realized by the network with the parameters  $\theta$ . The challenge lies in the  
740 intractability of the integral for complex networks.

741  
742 Sommerfeld et al. (2025) address this challenge by representing the BNN as a factor graph. Factor  
743 graphs are probabilistic graphical models that facilitate the approximation of marginal distributions.  
744 The predictive posterior, as an integral over a product of factors, is well suited for this representation.  
745 However, direct modeling of the neural network function with Dirac delta functions leads to intractable  
746 message computations within the factor graph. Therefore, Sommerfeld et al. (2025) adopt a scalar-  
747 level representation, introducing latent variables and elementary Dirac delta factors to model the  
748 neural network’s operations at a finer granularity. This decomposition allows for the derivation of  
749 MP equations, which are central to the belief propagation algorithm used to approximate marginals.

750  
751 Due to the complexity of exact message calculations in large networks and datasets, and the presence  
752 of cycles in the factor graph, Sommerfeld et al. (2025) employ several key approximations. Firstly,  
753 messages are approximated as scaled Gaussian densities, leveraging the property that the product of  
754 Gaussians results in another Gaussian. This is facilitated by the use of the natural parameterization  
755 of Gaussians, where the precision and precision-mean are used instead of the mean and variance.  
Specifically, for a Gaussian  $\mathcal{N}(\mu, \sigma^2)$ , the precision is defined as  $\rho = 1/\sigma^2$  and the precision-mean as

756  $\tau = \mu\rho$ . This representation simplifies the multiplication of Gaussian densities and vastly improves  
 757 the numerical stability of various message equations.

758 Secondly, moment matching is employed to approximate messages related to nonlinear activation  
 759 functions. Direct moment matching is used where feasible. In cases where direct moment matching is  
 760 difficult, a marginal approximation is used. This involves approximating the full marginal distribution  
 761 of a variable with a Gaussian, and then using this approximation to derive an approximate message.  
 762

763 Thirdly, variational message passing is used to approximate the messages associated with the product  
 764 of two variables. This approach, as described in Stern et al. (2009), helps in managing the symmetries  
 765 present in the posterior distribution of BNNs.

766 The training procedure in Sommerfeld et al. (2025) uses loopy belief propagation with a specific  
 767 message schedule to handle cyclic dependencies in the factor graph. A batching strategy is also  
 768 employed to manage the computational cost of large datasets. This involves processing the training  
 769 data in smaller batches and aggregating messages from inactive examples. Finally, prediction for  
 770 unseen inputs is performed by propagating messages from the training branches to the prediction  
 771 branch, effectively using the approximate posterior over weights as a prior during inference. This is  
 772 an advantage over IVON which can only produce samples of the posterior predictive distribution via  
 773 sampling from the posterior and doing a forward pass for each sample.

774 Lastly, the regression factor  $\mathcal{N}(a; y, \beta^2)$  is of relevance for this paper. It represents the likelihood of  
 775 observing a target value  $y$  given an input  $x$  and network parameters (or equivalently the networks  
 776 output). Its role is crucial in updating the posterior distribution of the parameters during training.  
 777 Formally, it connects two random variables  $a$ , a scalar outputted by the network, and  $y$ , the target. The  
 778 variance  $\beta^2$  is a constant (hyperparameter) we can interpret as aleatoric uncertainty or as a measure  
 779 of belief we have in our target, i.e., epistemic uncertainty.

## 781 C DETAILS ON THE IMPLEMENTATION

### 782 C.1 TRAINING IMPLEMENTATION

784 To test the impact of FG-BNNs, we developed an implementation of the concept above in Julia. The  
 785 choice for Julia was made because the factor graph framework is implemented in Julia, and this  
 786 programming language offers easy integration of common Python machine learning libraries. Our  
 787 `Bayesian A2C Agent` is compatible with `Gymnasium` environments, the successor of `OpenAI Gym`.  
 788 The current implementation only supports continuous action spaces.

789 In the main training loop, for each time step, the state of every environment is taken and used by the  
 790 actor to predict the action mean and standard deviation. From these values, the action is sampled and  
 791 executed, leading to new observations. Afterwards, the critic is run to estimate the values of the new  
 792 and previous state, for a vanilla advantage and return estimation. While many production actor-critic  
 793 implementations use the Generalized Advantage Estimator (GAE), we decided to go with the direct  
 794 way for simplicity. Afterwards, the pseudo targets are calculated as described Section 3.2, and the  
 795 parameter training happens for both actor and critic.

796 The implementation performs extensive logging of various metrics during training. After each batch,  
 797 24 measurements are written to `Weights & Biases` (Biewald, 2020). Examples are the outputs of the  
 798 actor and critic, the differences between predicted action means and pseudo targets, the precisions  
 799 and precision mean parameters from the inner workings of the FG-BNNs, and of course rewards. As  
 800 most of these are empirical distributions among one batch (like the taken actions), or the parameters  
 801 of the network, we save their mean, and 10, 25, 50, 75, and 90 percentiles. Rigorous logging also  
 802 helps to explain the findings and differences.

803 The factor graph itself is not part of the `Bayesian A2C Agent` class, but accessed via a standard  
 804 interface. For this interface, we needed to make slight changes to the factor graph library. The  
 805 adapted code is part of the repository that is linked in this paper. The only communication between  
 806 the agent and the factor graph happens through two function calls: `predict` returns the means and  
 807 standard deviation for a batch, and `train_batch_new` is called with a state and a target (where we give  
 808 the pseudo-target in case of the actor), together with the configuration parameters. This decoupling  
 809 between the algorithm and the factor graph makes it possible to easily switch the type of the neural  
 network to analyze the impact of the message-passing network compared to normal NNs and other

810 types of BNNs. In the experimental section, we will make use of this and run experiments with  
811 different kind of networks.

## 813 C.2 DESIGN DECISIONS FOR TRAINING STABILITY

814  
815 Just like for normal NNs, the samples in one batch should be independent and identically distributed  
816 in Bayesian networks. To satisfy this condition best, there are as many environments instantiated  
817 as the batch size says (in our case 256). Since many of our control problems have a fixed episode  
818 length, we perform a different number of pre-execution steps to each of the environments to make  
819 sure starts and ends are equally distributed among all batches. Additionally, in each step we reset an  
820 environment with a probability of 0.5% to avoid the formation of similar environments in one batch.  
821 If the environments were synchronized in terms of timing, the samples would not be stochastically  
822 independent.

823 Research in the field of actor critic algorithms has shown that a common problem is overshooting  
824 while training the actor. When the parameters are optimized according to the policy gradient theorem,  
825 the policy often changes “too much”, leading to unstable training. In response to that, algorithms  
826 like PPO aim to restrict the change of the policy. PPO clips the gradients to restrict the change of  
827 the action probability, slowing down the learning deliberately to improve stability. We adapt this  
828 idea by clipping the pseudo targets depending on the standard deviation given by the BNN. The  
829 hyperparameter  $\kappa$  defines this percentage, defaulting to 1.25.

830 In our analysis, we found that training stability could be improved by lowering the learning rate  
831 over the course of the training. Reducing the learning rate is a well-known technique in Machine  
832 Learning Smith (2018). Our algorithm uses an exponential decline, where the learning rate decreases  
833 exponentially with the number of steps, leveling off at 1/10 of the initial learning rate.

834 Apart from these described design decisions, the implementation does not use tunings to keep it  
835 simple and put the focus on the change in function approximators.

## 837 C.3 NOTABLE DIFFERENCES TO NORMAL NEURAL NETWORKS

838  
839 While many RL algorithms use the tanh-function to squeeze the output into the interval of  $[-1, 1]$ , we  
840 do not use tanh and go for clipping instead. This decision was made because a tanh activation function  
841 would squeeze the normal distribution heavily that comes as a result of the last layer. Especially if  
842 the values are far away from zero, tanh would squeeze the variance close to zero, resulting in a stop  
843 of exploration and change. Clipping to the allowed action space happens both during inference and  
844 during calculation of the pseudo targets, ensuring that the pseudo targets are not outside of the action  
845 space.

846 The architecture of the BNN is close to the architecture proposed in the PPO paper, with two hidden  
847 layers of 64 neurons each. The first layer’s input dimension is variable depending on the number  
848 of features of the selected environment, while the output of the critic has always one dimension,  
849 and the output of the actor has the same dimensionality as the action space. However, in contrast to  
850 the original implementation, the BNN uses LeakyReLU with a leak of 0.1 as its activation function.  
851 LeakyReLU was found to perform better in for FG-BNNs, as discussed by Sommerfeld et al. (2025).

## 852 C.4 HYPERPARAMETERS FOR FACTOR GRAPH BAYESIAN NEURAL NETWORKS (FG-BNNs)

853  
854 Furthermore, both the actor and the critic have a regression factor at their output position. This  
855 factor adds a default uncertainty to the prediction, resulting in two important hyperparameters  $\beta_{actor}$   
856 and  $\beta_{critic}$ . This hyperparameter sets a lower bound for the variance, which has the following two  
857 important consequences for the behavior of the algorithm:

- 858  
859 • When it comes to the actor, this is a minimum level of exploration, even if the rest of the  
860 algorithm has converged. It is important to not select the value too low to ensure that training  
861 does not get stuck too early.
- 862 • It controls the training speed of the algorithm. A low variance on the output level means a  
863 high prior towards the predicted value during the backward pass when training the parameters.  
If the prior is too prejudiced, adjustment of the parameters will be considerably delayed.

Especially in the case of the critic, this is a problem if a critic that is too self-confident and, therefore, slow to learn cannot keep pace with changes to the policy.

Apart from  $\beta_{actor}$  and  $\beta_{critic}$ , Actor Critic with BNNs requires some other choices of hyperparameters than standard deep learning. The choices made for our benchmarks are shown in Table 1. For PPO, we used the standard parameters given in Table 2.

Hyperparameter	Value	Description
<i>batch size</i>	256	The batch size for training equals the number of environments
<i>num batches</i>	1 000	It describes how many batches are trained before the training stops
$\gamma$	0.9	The trade-off parameter between immediate and long-term reward
$\alpha$	0.05	The “learning rate” in the Bayesian framework, effectively a scaling parameter for the advantage during pseudo-target calculation
<i>number factor</i>	3	It describes how often forward & backward messages are sent during the training of one batch
<i>graph iterations</i>		
<i>Preset</i>	0.005	A probability for resetting an environment at random
<i>action std expansion factor</i>	4.0	A regularization factor to foster exploration: the standard deviation of the actor gets multiplied by this before sampling an action during training
$\beta_{actor}$	0.05	Minimum uncertainty of the actor prediction
$\beta_{critic}$	0.4	Minimum uncertainty of the critic prediction
<i>action percentage clipping</i>	0.9	To improve stability, the pseudo targets are clipped to only 90% of the action space, ensuring that the side of the distribution also stays within the action space to a large percentage
$\kappa$	1.25	Clip factor to limit change of the actor in an update step
$\alpha_{PyTorch}$	0.003	The learning rate used in the PyTorch comparative implementation for both actor and critic

Table 1: Hyperparameters of the Bayesian Actor Critic (BAC) algorithm with their default values. In the bottom section, the hyperparameters for the PyTorch comparative algorithm are given.

Parameter	Value
Learning rate	$3 \cdot 10^{-4}$
Steps per update	2 048
Minibatch size	64
Epochs per update	10
<i>clip_range</i> ( $\epsilon$ )	0.2
Discount factor ( $\gamma$ )	0.9
Generalized advantage estimator factor ( $\lambda$ )	0.95
Entropy coefficient	0
Value function coefficient	0.5

Table 2: Hyperparameters for Proximal Policy Optimization (PPO).

## D ABLATION STUDY

The hyperparameters given in Table 1 needed to be chosen with care in order to achieve good performance. We give a few examples for the Pendulum environment to justify that our selection is reasonable. However, due to the large space of possible hyperparameters, we were unable to do a fully comprehensive analysis. It may be well possible that a different choice of hyperparameters will give better results. That is especially true for the variety of environments. We did not do environment-specific hyperparameter tuning to get as general results as possible. Hence, the environments may benefit from additional tuning.

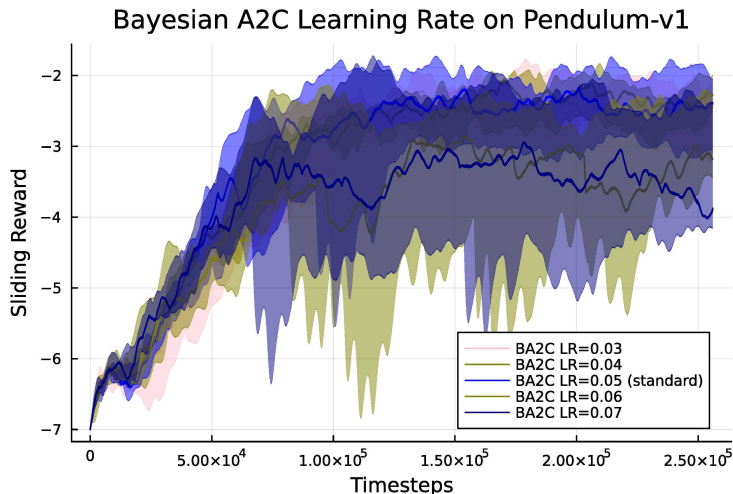


Figure 2: Learning curves for the BA2C algorithm with different learning rates on the Pendulum environment.

Figure 2 shows the learning curves of the BA2C algorithm on the Pendulum environment with different learning rates. The range of analyzed parameters is from 0.03 to 0.07, where 0.05 is our standard parameter. While these learning curves do not show completely different paths, there are noticeable differences. A lower learning rate tends to lead to slower learning, while the highest learning rate showed the highest initial training speed. However, the learning curves with a learning rate of 0.07 already tend to be more instable. This justifies our selection of 0.05 as our standard selection.

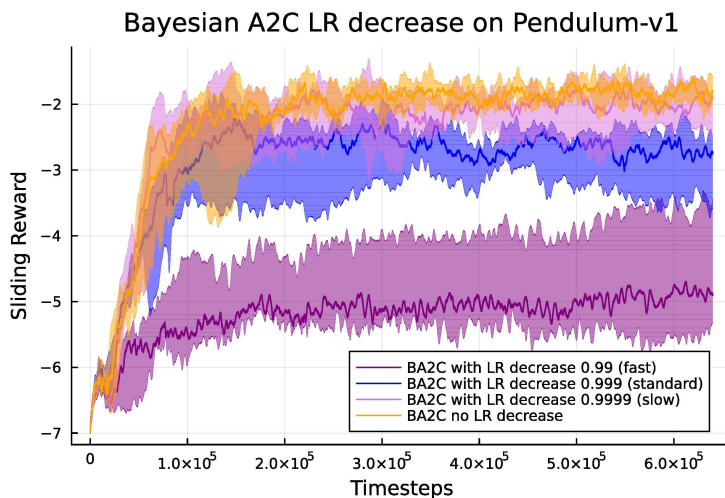


Figure 3: Learning curves of BA2C with different levels of decreasing learning rate on 640,000 steps.

Many machine learning implementations use a decreasing learning rate, so does our BA2C. We chose the exponential decline of the learning rate due to Figure 3. As you can see, the algorithm gets problems with stability if the learning rate remains unchanged, especially during later phases of training. Our standard parameter of 0.999 performs great in these experiments. If the learning rate is decreased faster, the initial learning rate will be too small, and no acceptable performance can be achieved.

When comparing the learning rate of the BA2C algorithm to the learning-rate hyperparameter in of neural networks trained with a classical gradient-based optimizer, please keep in mind that it has a very different semantic in BA2C. It is not a scaling factor for the gradients to directly influence the parameters, but it is the distance of the pseudo targets to the current prediction.

The following analysis of learning rates in the context of Adam or IVON neural networks refers to the learning rate given to the optimizer. The scaling factor on the distance of the pseudo targets stays at 0.05.

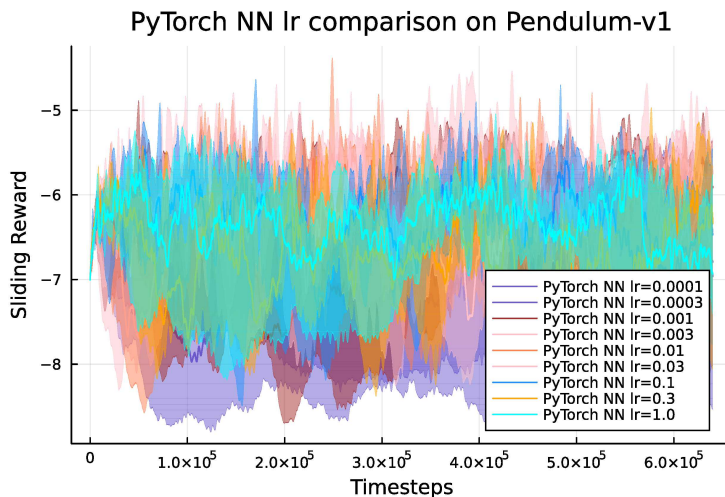


Figure 4: Comparison of different learning rates for PyTorch neural networks, selected for both actor and critic in our BA2C algorithm on 640,000 steps.

Next, we need to justify the learning rate of 0.003 for the PyTorch neural networks that are trained via Adam when they are used for actor and critic in our BA2C algorithm. First, 0.003 is also the chosen learning rate in Stable Baselines 3 for their neural networks of the same architecture, so we use an appropriate comparison here. Second, we checked the performance of the algorithm on the Pendulum environment while testing learning rates from 0.0001 to 1. The learning curves are visualized in Figure 4. As we can see, the chosen 0.003 clearly gives the best performance. To illustrate that it is actually able to reach a competitive performance, we show it for 640,000 steps.

According to the hyperparameter guide in the appendix of the IVON paper Shen et al. (2024), we needed to set the effective sample size. In our case, this is the number of training steps. Additionally, we needed to set the learning rate parameter. We chose 0.1, a common selection for small neural networks. As you can see in Figure 13, IVON was relatively effective for some of the environments. On Pendulum however, IVON did not work at all. To check if the selection of the learning rate was responsible for this, we ran IVON with the same variety of learning rates like in 4. Unfortunately, Figure 13 shows that no learning rate worked for IVON on Pendulum.

We decided to implement the algorithm with advantage normalization. The primary reason to do this is to avoid the need of environment-specific learning rates. We analyzed the impact of this normalization on Pendulum and found no significant difference, when the rest of the configuration stays the same. Figure 6 shows the learning curves.

The next hyperparameter to justify is  $\beta_{actor}$ , the default uncertainty the algorithm carries. The default choice is 0.05. We tested a few parameters (shown in Figure 7) higher and lower, and found that the standard selection works best. The results for  $\beta_{critic}$  are similar.

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

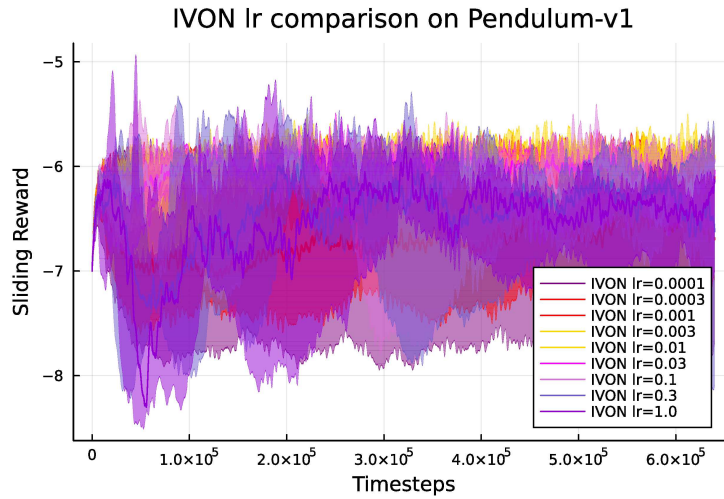


Figure 5: Comparison of different learning rates for BA2C training with IVON on 640,000 steps.

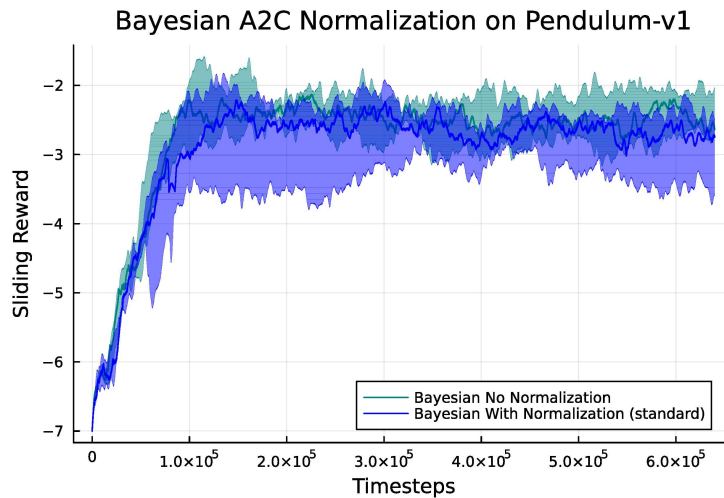


Figure 6: Learning curves of BA2C with and without advantage normalization.

To foster exploration, we use a linear scale between the uncertainty of the actor and the standard deviation of the stochastic, Gaussian policy. The standard choice of 4.0 is not backed by theoretic findings, but by empiric results. Figure 8 shows a few different selections of this parameter around 4.0.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

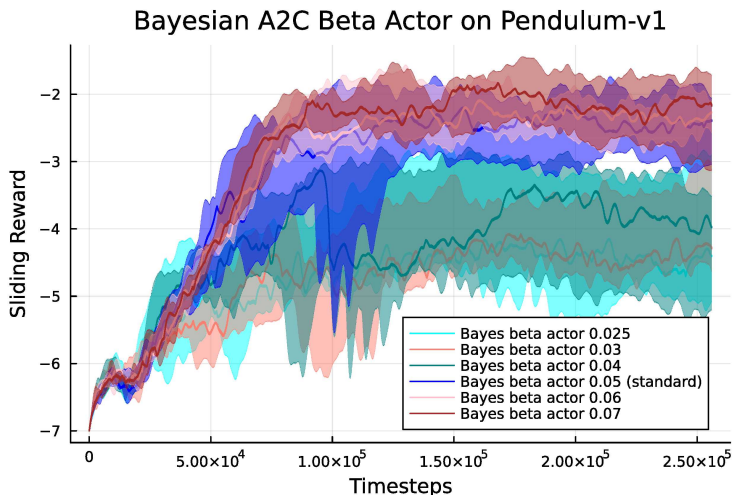


Figure 7: Different choices of  $\beta_{actor}$  result in these learning curves.

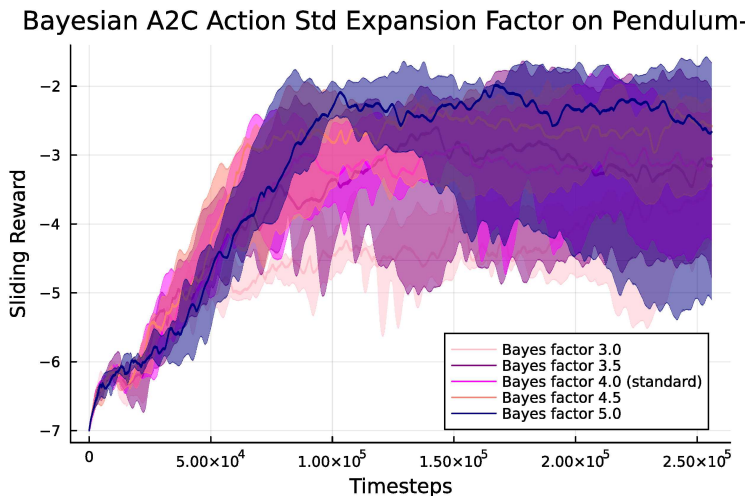


Figure 8: Learning curves of the BA2C algorithm with different scaling factors from the uncertainty to the standard deviation of the Gaussian policy.

## E ON THE FASTER LEARNING OF BA2C ON PENDULUM

Figure 9 shows the four combinations of FG-BNNs and PyTorch networks like Figure 1(a) in the main part. However, here we plot the curves over a longer time (640,000 samples) that illustrate the ability of the combination with two PyTorch networks to achieve competitive performance when more steps are granted.

For the same runs, Figure 10 shows the Mean Squared Error between the predicted state value of the critic and the target of the critic update. As we can see, this mean squared error is smaller by an order of magnitude for Bayesian critics (blue and yellow lines) compared to PyTorch critics (green and pastel-colored lines).

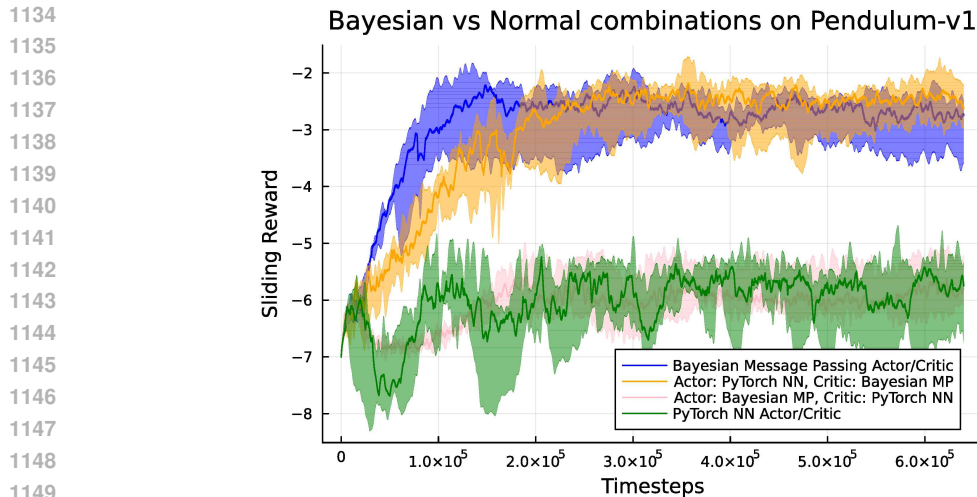


Figure 9: Learning Curves of BA2C combinations on Pendulum on a horizon of 640,000 steps.

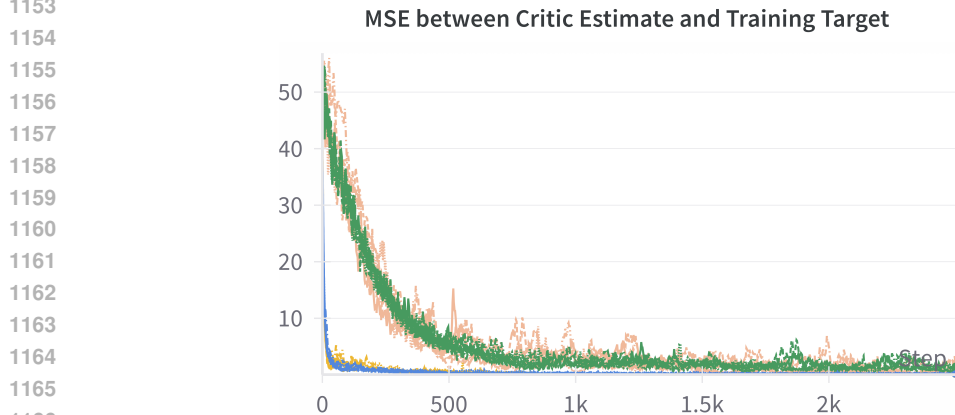


Figure 10: MSE between predicted state value by the critic and target: In the combinations with a Bayesian critic, the MSE decreases significantly faster, which is analogous to the finding in Figure 9, where agents with Bayesian critic achieved higher performance much quicker. Colors are also chosen like in Figure 9.

## F ADDITIONAL EVALUATION ON OTHER GYMNASIUM ENVIRONMENTS

We continue the evaluation from Section 5 here. In the main paper, we analyzed the Pendulum environment with respect to different aspects. We compared FG-BNNs against PyTorch neural networks within the BA2C algorithm (Section 5.1), verified the importance of uncertainty-driven exploration (Section 5.2), tested the FG-BNNs against IVON (Section 5.3), and finally tested BA2C against popular PPO implementations. Here, we will add plots for the same experiments on 8 total environments (Pendulum plus seven other Gymnasium environments). Figure 11 shows how the four combinations of FG-BNNs and PyTorch neural networks perform on the eight environments. Figure 12 shows the comparison between networks with uncertainty-driven exploration and without. The comparison to IVON for the eight environments is given in Figure 13. In Figure 14, we provide the comparison of BA2C against the popular libraries.

While the plots are organized by experiment, the following text will walk through the results by environment. We will highlight and explain notable insights.

- (a) The Pendulum environment has been extensively studied in the main part.

- 1188 (b) The Mountain Car environment is special because it is a setup with a hidden reward. The  
 1189 agent has to swing a car back and forth several times to climb a mountain. The reward  
 1190 is only given when the mountain has been successfully climbed, and each acceleration  
 1191 reduces the reward. An optimal agent would swing the car up the hill and accept the cost.  
 1192 However, Mountain Car is an environment that usually requires memorization (storing old  
 1193 experience). This would help the agent to learn from even very rare successes. Sometimes,  
 1194 we see upward-spikes in the learning curves for some of the agents. These upward-spikes  
 1195 arise if the environment was solved and a high reward was granted. However, all of our  
 1196 candidates are on-policy algorithms without memorization. Usually, they cannot learn from  
 1197 these rare successes, so they optimize to a policy with 0 actions. Our BA2C algorithms all  
 1198 find out about this very quickly. There are only marginal differences between the different  
 1199 networks and options within BA2C. When compared to the libraries, BA2C is the fastest  
 1200 algorithm to discovering the 0-policy, and achieves this policy practically instantaneous.  
 1201 Dopamine and SB3 take significantly longer. The reason for the lower final-performance  
 1202 of the Bayesian A2C is its forced level of higher exploration, causing costs for the taken  
 1203 actions.
- 1204 (c) The Bipedal Walker environment is about moving a two-dimensional robot in a specific  
 1205 direction while controlling four joints. On this task, all libraries and our BA2C heavily  
 1206 struggle with stability, as it can be seen in Figure 14(c). This is not a surprise, given the  
 1207 24-dimensional observation space and four-dimensional action space. An average reward of  
 1208 0 means the agent can stabilize without falling, but does not move forward. This is a level  
 1209 that is only temporarily, never reliably achieved by all implementations. On the horizon of  
 1210 50,000 steps, there is no clear winner. (On a much longer horizon, SB3 is able to improve  
 1211 significantly).
- 1212 (d) The Inverted Pendulum environment can be considered the easiest out of these. BA2C  
 1213 achieves basically optimal performance, with few runs suffering from minor stability prob-  
 1214 lems. Within the different options of implementing BA2C, we can confirm our findings  
 1215 from Pendulum: FG-BNNs work best, missing out on uncertainty-driven exploration hurts  
 1216 stability, and IVON gives acceptable but non-optimal performance. SB3 and Ray RLlib give  
 1217 high performance. SB3 is slightly behind in the beginning but achieves optimal performance  
 1218 most reliably, while Ray RLlib is fast in the beginning but not all runs get optimal. Dopamine  
 1219 also improves performance but does not get to a competitive level.
- 1220 (e) We see a picture with very similar findings when testing Inverted Double Pendulum. One  
 1221 noticeable difference is in the comparison to the libraries (Figure 14(e)). Here, RLlib is a  
 1222 touch faster than BA2C in the very beginning, but BA2C outperforms RLlib after around  
 1223 4,000 samples. Stable Baselines 3 is significantly slower, but achieves significantly better  
 1224 performance after 12,500 steps. One BA2C run became instable in a later phasis.
- 1225 (f) On the Reacher environment, BA2C is clearly faster than all PPO libraries (Figure 14(f)).  
 1226 However, it is interesting to see that the choice of FG-BNNs for both actor and critic is  
 1227 not optimal within the options for BA2C. As Figure 11(f) shows, the selection of PyTorch  
 1228 neural networks for the actor results in better performance, and so does IVON (Figure 11(f)).  
 1229 The reason might be the problem of FG-BNNs to handle multi-dimensional outputs well.  
 1230 Nevertheless, BA2C is a highly effective algorithm with one of the other options selected.
- 1231 (g) The plot for Swimmer (Figure 14(g)) shows a remarkable hump at relatively early training  
 1232 for all variants of BA2C. We investigated the reason for this, and found that it is caused  
 1233 by the way the learning curves are designed, not by actual policy changes. The reason is  
 1234 that our BA2C runs with 256 environments in parallel. They are at a similar initialization-  
 1235 like state in the beginning. Then, the agent quickly finds out a movement that brings the  
 1236 swimmer forward by a small distance by contracting the arms. This is executed for all 256  
 1237 environments. However, when the Swimmer is in the contracted state, it stays in this state,  
 1238 yielding zero rewards from now on. Unlike the other environments, Swimmer and Half  
 1239 Cheetah have very long episodes (a few thousand steps). The 256,000 steps that are shown  
 1240 in the learning curve however, are the result of round-robin scheduling of our environments,  
 1241 where only 1,000 steps are executed in each environment. Hence, BA2C’s learning curve  
 draws a pattern like a single evaluation on an environment. Why is this not a problem for  
 the other environments? First, the episodes on the other environments are much shorter.  
 Second, the pre-execution takes care of bringing environments with different progress into

1242 one batch. But random sampling does not change the initial state of Swimmer much, while  
1243 also running too few samples for bringing it close to its terminal state. Third, are less prone  
1244 to entering deadlocks. For the stabilization tasks, the episode is often terminated early, and  
1245 not stuck in an unhelpful state if it was entered. However, these aspects just impact the  
1246 rendering of the plots, not the actual workings of the algorithms.

1247 (h) On Half Cheetah, we see a similar hump like on Swimmer, and we see Half Cheetah being  
1248 outperformed by the libraries Stable Baselines 3 and RLlib (Figure 14(h)). IVON and normal  
1249 neural networks also do not work (Figure 13(h)), but we see remarkable performance of the  
1250 combination of a Bayesian critic with a PyTorch actor in Figure 11(h) which is competitive  
1251 to Stable Baselines 3. This indicates that there is indeed a huge potential if the problem of  
1252 numeric instability of the factor-graph framework is solved. Moreover, this highlights the  
1253 performance of our BA2C algorithm even on more complex environments.

1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

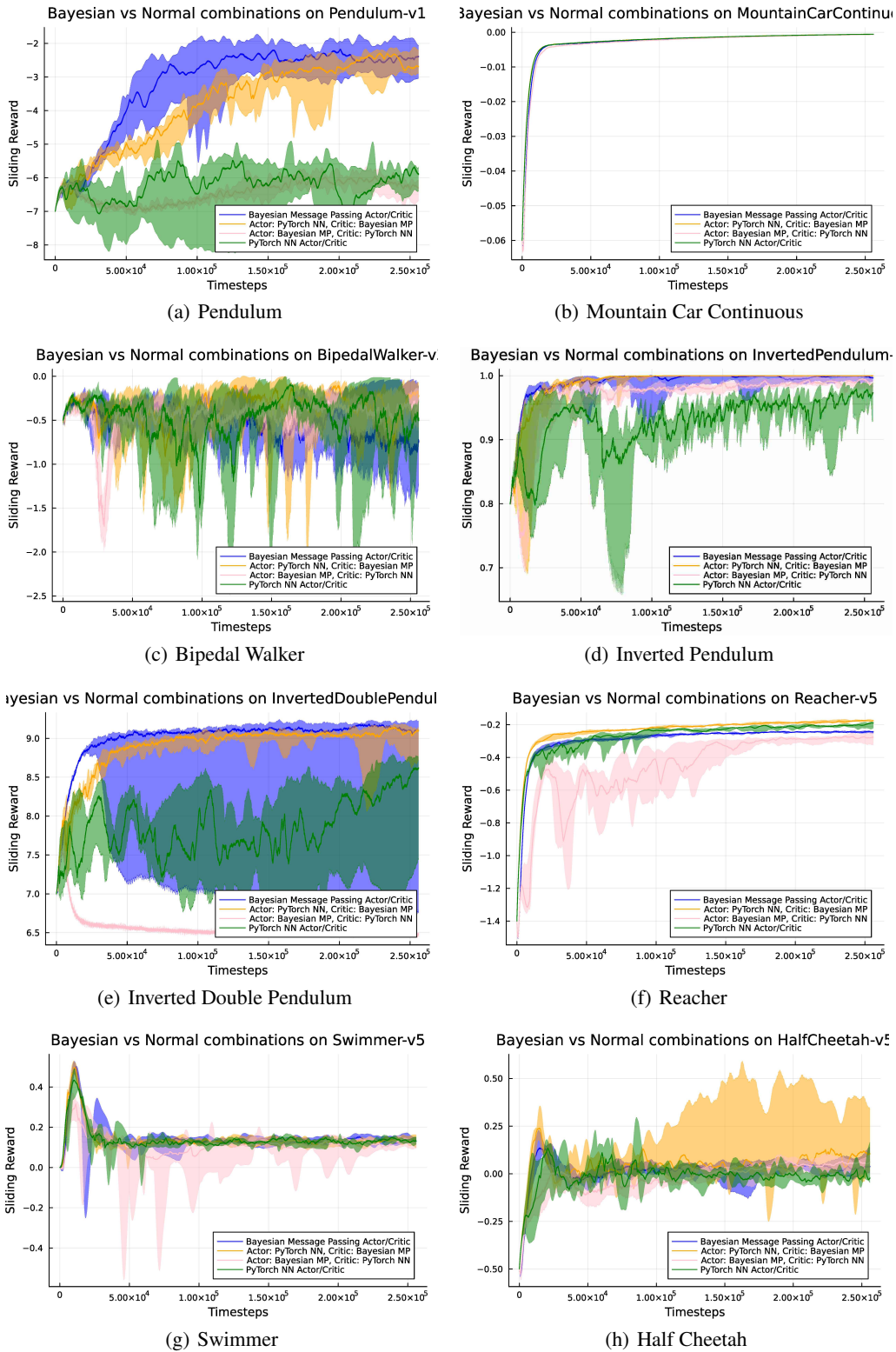


Figure 11: Combinations of Factor Graph & PyTorch.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

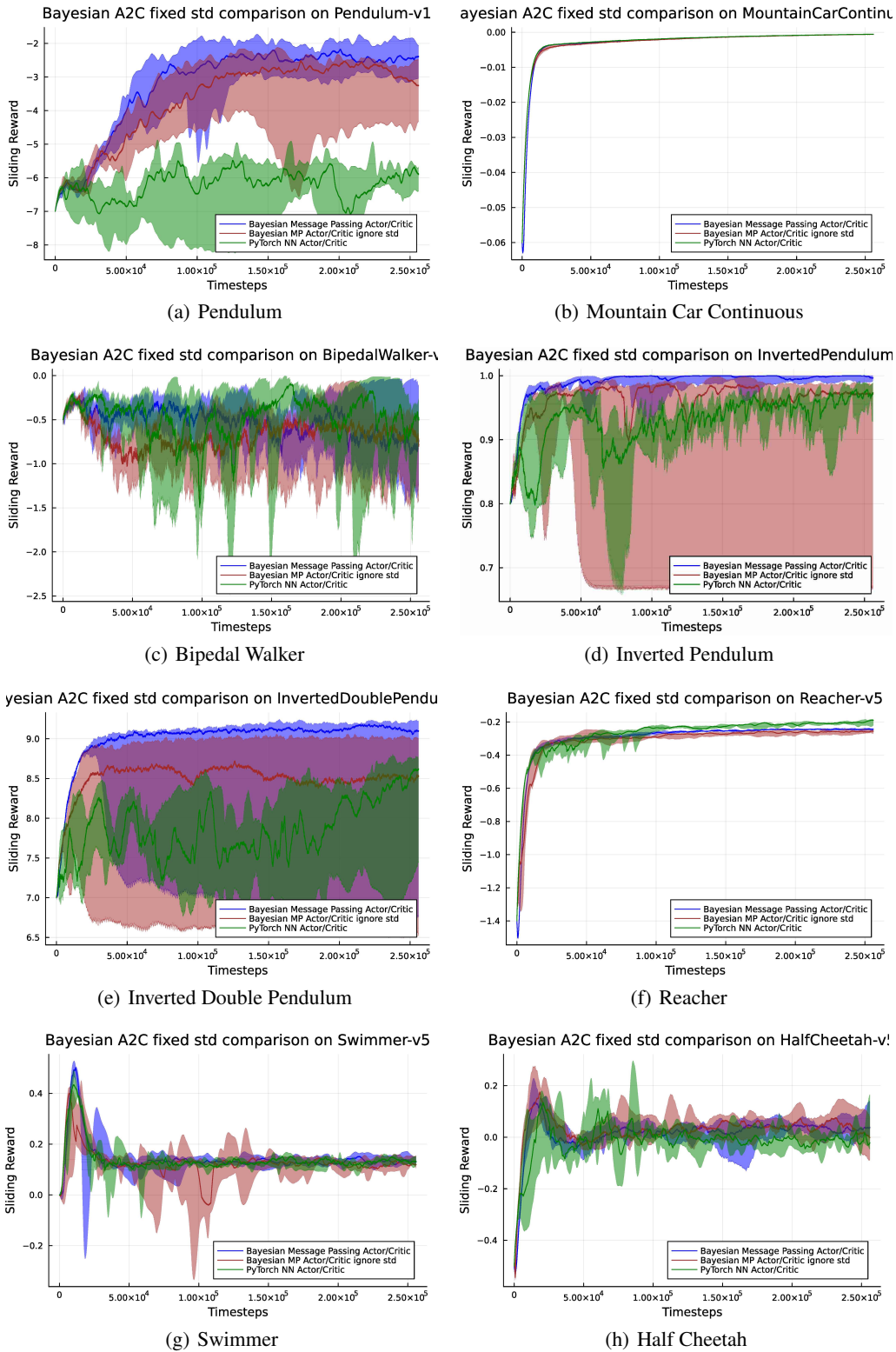


Figure 12: Factor Graph with Uncertainty-based Exploration vs. Fixed Exploration.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

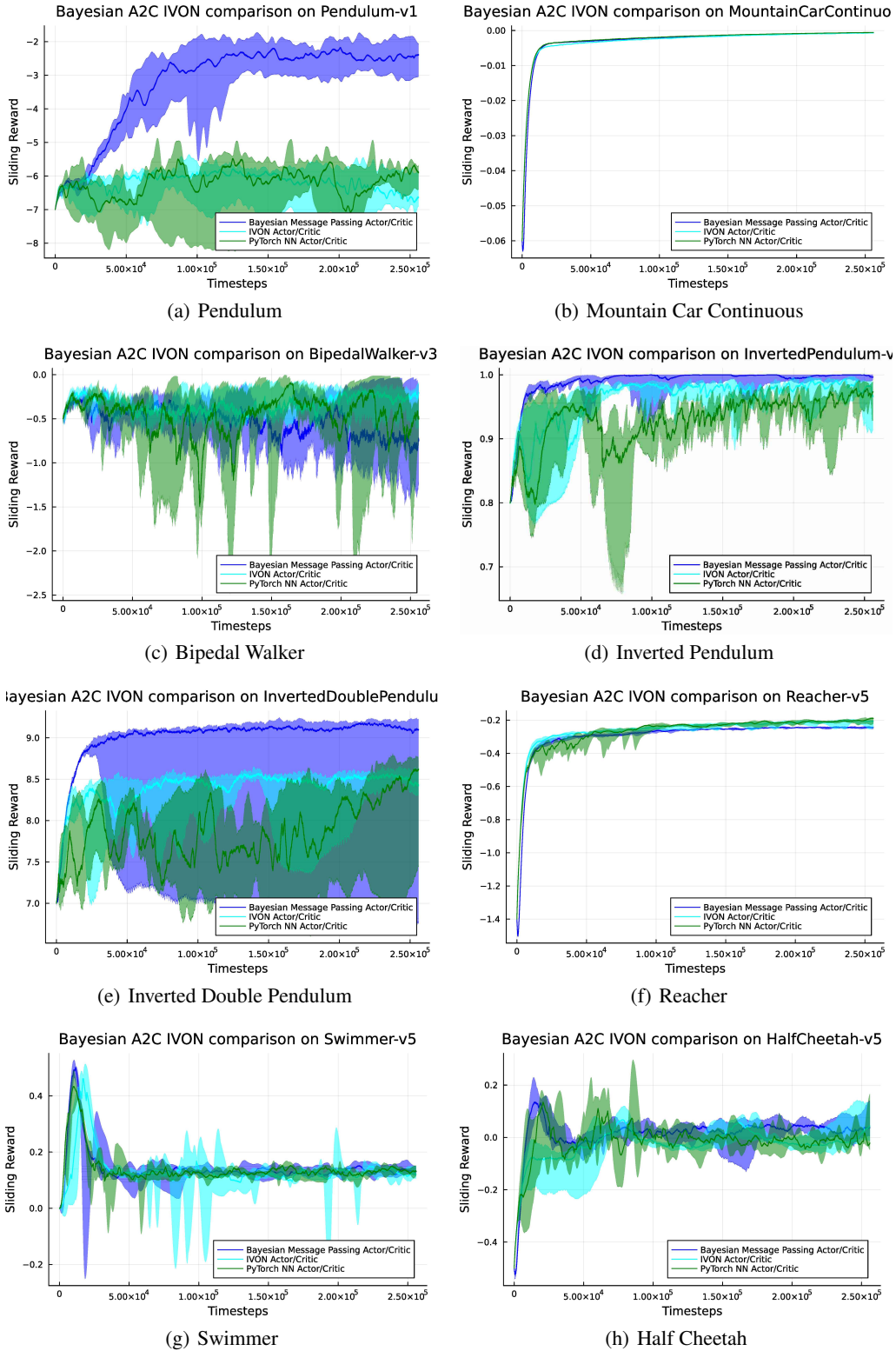


Figure 13: Learning curves of the eight environments to compare IVON against our standard BA2C and PyTorch neural networks.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

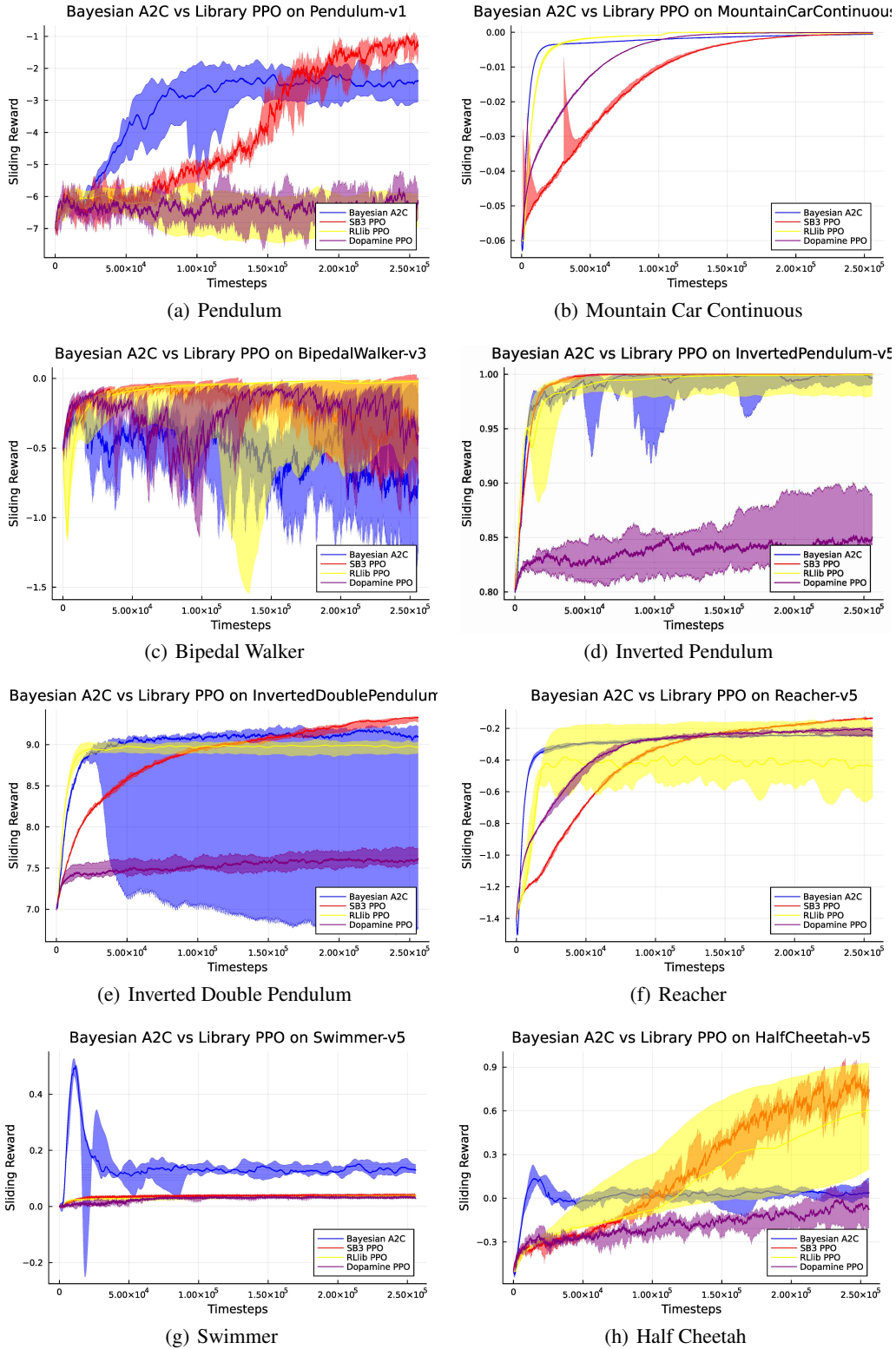


Figure 14: BA2C (blue) against PPO agents from three common libraries: SB3 (red), Ray RLlib (yellow), and Dopamine (purple).