
Disentangling Epistemic and Aleatoric Uncertainty in Reinforcement Learning

Bertrand Charpentier¹, Ransalu Senanayake², Mykel Kochenderfer², Stephan Günnemann¹

Department of Informatics & Munich Data Science Institute, Technical University of Munich¹

Stanford Intelligent Systems Laboratory, Stanford University²

{charpent, guennemann}@in.tum.de, {ransalu, mykel}@stanford.edu

Abstract

Characterizing *aleatoric* and *epistemic* uncertainty on the predicted rewards can help in building reliable reinforcement learning (RL) systems. Aleatoric uncertainty results from the irreducible environment stochasticity leading to inherently risky states and actions. Epistemic uncertainty results from the limited information accumulated during learning to make informed decisions. Characterizing aleatoric and epistemic uncertainty can be used to speed up learning in a training environment, improve generalization to similar testing environments, and flag unfamiliar behavior in anomalous testing environments. In this work, we introduce a framework for disentangling aleatoric and epistemic uncertainty in RL. **(1)** We first define four *desiderata* that capture the desired behavior for aleatoric and epistemic uncertainty estimation in RL at both training and testing time. **(2)** We then present four RL *models* inspired by supervised learning (i.e., Monte Carlo dropout, ensemble, deep kernel learning models, and evidential networks) to instantiate aleatoric and epistemic uncertainty. Finally, **(3)** we propose a practical *evaluation* method to evaluate uncertainty estimation in model-free RL based on detection of out-of-distribution environments and generalization to perturbed environments. We present *theoretical* and *experimental* evidence to validate that carefully equipping model-free RL agents with supervised learning uncertainty methods can fulfill our desiderata.

1 Introduction

An agent is expected to satisfy three important properties for a reliable deployment in real-world applications: **(i)** The agent should learn *fast* with as few episode failures as possible. **(ii)** The agent should *maintain high reward* when facing new environments similar to the training environment after deployment. **(iii)** The agent should *flag anomalous environment states* when it does not know what action to take in an unknown environment. These three practically desirable properties translate into three technical properties in reinforcement learning agents. Indeed, a reinforcement learning agent should achieve high *sample efficiency* at training time [1], high *generalization* performance on test environments similar to the training environment [2], and high *Out-Of-Distribution (OOD) detection* scores on environment unrelated to the training task [3], [4].

In this paper, we argue that *aleatoric* and *epistemic* uncertainty are key concepts to achieve these desired practical and technical properties. The aleatoric uncertainty represents the irreducible and inherent stochasticity of the environment. Thus, an environment region with high aleatoric uncertainty is unlikely to be interesting to explore at training time because it could be uninformative (e.g. a sensor is very noisy) or dangerous (e.g. the environment has an unpredictable behavior). In contrast, the epistemic uncertainty represents the lack of information for accurate prediction. Thus, an environment region with high epistemic uncertainty is potentially promising to explore to build a better understanding of the environment (e.g., a state has unknown transition dynamics because it has never been explored).

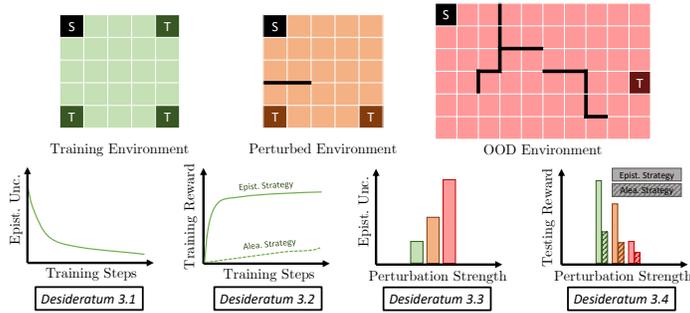


Figure 1: Overview of our proposed desiderata for uncertainty in RL (See sec. 3).

The core motivation of our work is to disentangle the properties of aleatoric and epistemic uncertainty estimates in RL to build agents with reliable performance in real-world applications. This motivation is similar to supervised learning (SL) where previous works defined *desiderata*, *models*, and *evaluation* methods for aleatoric and epistemic uncertainty [5]–[8]. Important examples of models using a single or multiple forward passes for uncertainty estimation in SL are MC dropout [9], ensemble [10]–[12], deep kernel learning [13]–[16], and evidential networks [17]–[20]. Further, empirical evaluation of uncertainty estimates in SL focuses *only* on testing time with Out-Of-Distribution (OOD) detection and generalization or detection of shifts [7], [21]. In contrast to SL, the RL setting is more complex since it cares about the performance of uncertainty estimates at *both* training and testing time.

Our Contributions. In this work, we propose a framework for aleatoric and epistemic uncertainty estimation in RL: **(Desiderata)** We explicitly define four desiderata for uncertainty estimation in RL at both *training* and *testing time* (See fig. 1). They cover the behavior of *aleatoric* and *epistemic* uncertainty estimates w.r.t. the sample efficiency in the training environment and, w.r.t. the generalization performance in different testing environments. **(Models)** We carefully combine a diverse set of uncertainty estimation methods in SL (i.e. MC dropout, ensemble, deep kernel learning, and evidential networks) with Deep Q-Networks (DQN) [22], a ubiquitous RL model that is not equipped with uncertainty estimate by default. These combinations require a *minimal* modification to the training procedure of the RL agent. We discuss *theoretical* evidence on the ability of these combinations to fulfill the uncertainty desiderata. **(Evaluation)** Finally, we also propose a *practical* methodology to evaluate uncertainty in RL based on OOD environments and domain shifts.

2 Problem Setup

Uncertainty in SL. The objective of SL is to accurately predict the output $y^{(i)}$ given an input $\mathbf{x}^{(i)}$ with index i . It differentiates between two types of uncertainty: the uncertainty on the label prediction $y^{(i)}$ described by the *aleatoric distribution* $\mathbb{P}(y^{(i)} | \boldsymbol{\theta}^{(i)})$ with parameters $\boldsymbol{\theta}^{(i)}$ estimated from the input $\mathbf{x}^{(i)}$, and the uncertainty on the predicted label distribution parameters $\boldsymbol{\theta}^{(i)}$ described by the *epistemic distribution* $\mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{(i)})$ with parameters $\boldsymbol{\chi}^{(i)}$ estimated from the input $\mathbf{x}^{(i)}$. Intuitively, the variations of the aleatoric distribution will be high when the input $\mathbf{x}^{(i)}$ does not provide discriminative information to determine the label $y^{(i)}$. The variations of the epistemic distribution will be high when the input $\mathbf{x}^{(i)}$ does not provide enough information to determine the label distribution $\mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{(i)})$ described by the parameters $\boldsymbol{\chi}^{(i)}$. Thus, the aleatoric uncertainty can be measured in practice by computing the entropy, i.e. $u_{\text{alea}}(\mathbf{x}^{(i)}) = \mathbb{H}(\mathbb{P}(y^{(i)} | \boldsymbol{\theta}^{(i)}))$, or the variance, i.e. $u_{\text{alea}}(\mathbf{x}^{(i)}) = \text{Var}(\mathbb{P}(y^{(i)} | \boldsymbol{\theta}^{(i)}))$, of the aleatoric distribution, while the epistemic uncertainty can be measured by computing the entropy, i.e. $u_{\text{epist}}(\mathbf{x}^{(i)}) = \mathbb{H}(\mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{(i)}))$, of the epistemic distribution [17]–[19]. Further, SL also distinguishes between sampling-based which often require multiple forward passes for uncertainty estimation, and sampling-free methods which often require a single forward pass for uncertainty estimation. Sampling-based methods like MC drop-out [9] and ensemble [10]–[12] estimate uncertainty by aggregating statistics (e.g. mean and variance) from different samples which *implicitly* describe the epistemic distribution $\mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{(i)})$. Sampling-free methods like deep kernel learning models [13]–[16] and evidential networks [17]–[20] estimate uncertainty by *explicitly* parametrizing the epistemic distributions $\mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{(i)})$ with known distributions such as Normal and Normal Inverse-Gamma (NIG) distributions, thus enabling efficient and closed-form computation of the distribution statistics (incl. mean, variance or entropy).

Uncertainty in RL. We consider the task of learning RL policies interacting with an environment with fully observed states at every time step t . The environment is described by a Markov Decision Process (MDP) $(S, A, r, T, \rho, \gamma)$ where S is the state space, A is the action space, $r(s^{(t)}, a^{(t)})$ is the

reward associated to the action $a^{(t)}$ and state $s^{(t)}$, $T(s^{(t+1)}|s^{(t)}, a^{(t)})$ is the transition probability, $\rho(s^{(0)})$ is the initial state distribution, and γ is the discount factor. Given the current state $s^{(t)}$, our goal is to learn a policy predicting the action $a^{(t)}$ leading to the highest reward $y^{(t)} = r(s^{(t)}, a^{(t)})$ in addition to the aleatoric uncertainty $u_{\text{alea}}(s^{(t)}, a^{(t)})$ and the epistemic uncertainty $u_{\text{epist}}(s^{(t)}, a^{(t)})$ on the predicted reward. Similar to SL, the aleatoric and epistemic distributions can be instantiated with $\mathbb{P}(y^{(t)} | \theta^{(t)})$ and $\mathbb{Q}(\theta^{(t)} | \chi^{(t)})$ where the predicted value is the future reward i.e. $y^{(t)} = r(s^{(t)}, a^{(t)})$. Intuitively, the variation of the aleatoric distribution will be high when the current state $s^{(t)}$ and action $a^{(t)}$ only contains noisy information to determine the future reward $y^{(t)} = r(s^{(t)}, a^{(t)})$ while the epistemic uncertainty will be high when the current state $s^{(t)}$ and action $a^{(t)}$ does not provide enough information according to the model to determine the reward distribution $\mathbb{Q}(\theta^{(t)} | \chi^{(t)})$ described by the parameters $\chi^{(t)}$. Finally, we consider three action selection strategies which is a crucial choice for exploration and generalization in RL: the *epsilon-greedy* strategy [23] which selects the action with the highest predicted reward with probability $1 - \epsilon$ and samples a random action otherwise, the *sampling-aleatoric* strategy which takes the action with the highest predicted reward based on one aleatoric distribution sample i.e. $a^{(t)} = \max_a y^{(t)}$ where $y^{(t)} \sim \mathbb{P}(y^{(t)} | \theta^{(t)})$, or the *sampling-epistemic* strategy which takes the action with the highest predicted reward based on one epistemic distribution sample i.e. $a^{(t)} = \max_a \mathbb{E}_{\mathbb{P}(y^{(t)} | \theta^{(t)})}[y^{(t)}]$ where $\theta^{(t)} \sim \mathbb{Q}(\theta^{(t)} | \chi^{(t)})$. The sampling-epistemic strategy corresponds to the Thompson sampling strategy [24].

3 Desiderata for Uncertainty Quantification in RL

In this section, we *explicitly* define four intuitive and general desiderata that capture the desired behavior for uncertainty estimates in our RL setup. The desiderata cover *aleatoric* and *epistemic* uncertainty at both *training* and *testing* time. The first distinction differentiates between aleatoric and epistemic uncertainty which are commonly used concepts in SL [5], [18], [19]. In contrast, RL mostly focuses on measuring aleatoric uncertainty with risk-sensitive policy or distributional RL [25]–[27]. The second distinction differentiates between training and testing time relevant to sample efficiency and generalization in RL. In contrast, SL mostly focuses on testing time performance.

Training Time. We describe the desired behavior of uncertainty estimates at training time. First, we describe the desired uncertainty behavior when observing more samples of the training environment.

Desideratum 3.1. *An agent training longer on states sampled from one specific environment should become more epistemically confident when predicting actions on states sampled from the same specific environment.*

Intuitively, an agent observing more samples from the same environment distribution should accumulate more knowledge through time thus being more epistemically certain. In practice, des. 3.1 expresses that the epistemic uncertainty estimates should reflect the accumulated knowledge, and thus the convergence, of the agent during training. We test des. 3.1 in the experiments (see sec. 5) by tracking the epistemic uncertainty at training time. Second, we describe the behavior of the total reward when selecting actions based on uncertainty estimates at training time.

Desideratum 3.2. *All else being equal, an agent selecting actions with the sampling-aleatoric strategy at training time should achieve lower sample efficiency than an agent selecting actions with the sampling-epistemic strategy.*

Intuitively, an agent exploring states with more (irreducible) aleatoric uncertainty would gain less knowledge about the environment dynamic than an agent exploring states with high epistemic uncertainty where the agent lacks knowledge. However, there is an important trade-off between over- or under-exploring epistemically uncertain actions which could lead to lower sample efficiency. The sampling-epistemic strategy, which corresponds to Thompson sampling [24], mitigates this exploration-exploitation problem by sampling action w.r.t. the epistemic distribution. Thompson sampling has already *empirically* demonstrated high sample efficiency in deep RL problems [9] and provably achieve low regret in many decision-making problems like multi-arms Bandit [28], [29]. In practice, des. 3.2 suggests that an agent should use the sampling-epistemic strategy for a better exploration-exploitation trade-off. We test des. 3.2 in the experiments (see sec. 5) by comparing the sample efficiency of the sampling-aleatoric and the sampling-epistemic strategies during training.

Testing Time. We describe the desired behavior of uncertainty estimates at testing time. First, we describe the desired uncertainty behavior when observing samples from an environment different from the training environment.

Desideratum 3.3. *At testing time, epistemic uncertainty should be greater in environments that are very different from the original training environments.*

The environment difference could be measured with different distances depending on the task or application requirements [30]. Intuitively, an agent should be less confident when observing new states at test time that were not used to collect knowledge at training time. In practice, des. 3.3 suggests that an agent should be able to use epistemic uncertainty estimates to detect states which are abnormal compared to the states observed during training. We test des. 3.3 in the experiments (see sec. 5) by comparing the epistemic uncertainty of the training environment against the uncertainty in noisy environments at testing time. Noisy environments include environments with completely random states and thus irrelevant to the training task, and environments with different strengths of perturbation on the original states, actions, or transition dynamics. Second, we describe the behavior of the total reward when selecting at testing time actions based on uncertainty aleatoric or epistemic uncertainty estimates.

Desideratum 3.4. *All else being equal, an agent sampling actions from the epistemic uncertainty at training and testing time should generalize better at testing time than an agent sampling actions from the aleatoric uncertainty.*

Intuitively, an agent exploring more epistemically uncertain states at training time would collect more knowledge about the environment, thus generalizing to more states at testing time. Further, since the environment dynamic is not directly observed, an agent should account for the epistemic uncertainty on the current state to take actions that generalize better at testing time. In particular, it has been shown that the Bayes-optimal Markovian policy at testing time is stochastic in general due to the partially observed MDP dynamic sometimes called epistemic POMDP [2]. In practice, des. 3.4 suggests that an agent should use the sampling-epistemic strategy for more robust generalization performance. We test des. 3.4 in the experiments (see sec. 5) by comparing the reward obtained by the sampling-aleatoric and the sampling-epistemic strategies at testing time. The two latter desiderata 3.3 and 3.4 express an important trade-off between assigning high uncertainty and generalizing to new test environments. Since an agent cannot generalize to all new environment because of the No Free Lunch Theorem [31], an agent should assign higher uncertainty to environments where it does not generalize. We jointly test des. 3.3 and des. 3.4 in the experiments (see sec. 5) by tracking the reward and the uncertainty estimates in test environments with different perturbation strengths.

4 Models for Uncertainty Quantification in RL

Model-free RL agents commonly rely on learning the expected return $Q^\pi(\mathbf{s}^{(t)}, a^{(t)})$ associated with taking action $a^{(t)}$ in state $\mathbf{s}^{(t)}$ and then following a policy π . It is defined by the Bellman equation: $Q^\pi(\mathbf{s}^{(t)}, a^{(t)}) = r(\mathbf{s}^{(t)}, a^{(t)}) + \gamma \mathbb{E}_{T, \pi}[Q^\pi(\mathbf{s}^{(t+1)}, a^{(t+1)})]$. Similarly, the optimal policy π^* achieving the highest expected reward satisfies the optimal Bellman equation [32]:

$$Q^{\pi^*}(\mathbf{s}^{(t)}, a^{(t)}) = r(\mathbf{s}^{(t)}, a^{(t)}) + \gamma \mathbb{E}_T[\max_{a^{(t+1)}} Q^{\pi^*}(\mathbf{s}^{(t+1)}, a^{(t+1)})] \quad (1)$$

However, the exact computation of the optimal Q -value is often intractable for large action or state spaces. Therefore, deep RL agents like DQN [22], PPO [33] and A2C [34] aim at approximating the optimal Q -value $Q^{\pi^*}(\mathbf{s}^{(t)}, a^{(t)})$ with a neural network $f_\theta(\mathbf{s}^{(t)}, a^{(t)})$ with parameter θ . In particular, DQN enforces eq. 1 by minimizing the squared temporal difference (TD) error $\|r(\mathbf{s}^{(t)}, a^{(t)}) + \gamma \max_{a^{(t+1)}} f_{\theta'}(\mathbf{s}^{(t+1)}, a^{(t+1)}) - f_\theta(\mathbf{s}^{(t)}, a^{(t)})\|_2$, where f_θ is the learned prediction network and $f_{\theta'}$ is the frozen target network regularly updated with the prediction network parameters during training. The TD error minimization is similar to SL regression with a MSE loss between the prediction $\hat{y}^{(t)} = f_\theta(\mathbf{s}^{(t)}, a^{(t)})$ and the target $y^{(t)} = r(\mathbf{s}^{(t)}, a^{(t)}) + \gamma f_{\theta'}(\mathbf{s}^{(t+1)}, a^{(t+1)})$ with the key difference that the exploration strategy select the targets that will be used during training.

Model-free Deep RL agents often show important limitations for uncertainty estimation because of their neural network architecture choice. For, instance, while DQN only outputs a single scalar representing the mean Q -value with no uncertainty estimates, PPO and A2C policies parameterized with standard ReLU networks would provably produce overconfident predictions for extreme input states [35]. In this work, we focus on equipping the widely used DQN RL agent with reliable uncertainty estimates. To this end, we combine DQN with four SL architectures for uncertainty estimation (incl. MC dropout, ensemble, deep kernel learning, and evidential networks) covering a diverse range of sampling-based and sampling-free methods. These four DQN combinations allow to instantiate both aleatoric and epistemic uncertainty with minimal modifications to the training procedure. We provide a summary of the uncertainty properties of these models in Tab. 1.

Table 1: Summary of the uncertainty properties of the models.

	DropOut	Ensemble	Deep Kernel Learning	Evidential Networks
Uncertainty concentration (Des. 3.1)	✗	✗	✗	✓
Alea. vs epist. sampling at training time (Des. 3.2)	✓	✓	✗	✓
OOD detection (Des. 3.3)	✗	✗	✓	✓
Alea. vs epist. sampling at testing time (Des. 3.2)	✓	✓	✗	✓

MC Dropout. DQN is combined with MC Dropout [9] in three steps: **(1)** it samples K independent set of model parameters $\{\theta_k\}_{k=1}^K$ by dropping activations with probability p , **(2)** it performs K forward passes $\mu_k, \sigma_k = f_{\theta_k}(s^{(t)}, a^{(t)})$, and **(3)** it aggregates predictions to form the mean prediction $\mu(s^{(t)}, a^{(t)}) = \frac{1}{K} \sum_{k=1}^K \mu_k$, the aleatoric uncertainty estimate $u_{\text{alea}}(s^{(t)}, a^{(t)}) = \frac{1}{K} \sum_{k=1}^K \sigma_k$, and the epistemic uncertainty estimate $u_{\text{epist}}(s_t, a_t) = \frac{1}{K} \sum_{k=1}^K (\mu_k - \mu(s^{(t)}, a^{(t)}))^2$. In this case, the aleatoric distribution is Gaussian while the epistemic distribution is implicitly represented by the sampled parameters $\{\theta_k\}_{k=1}^K$. Further, the sampling-epistemic strategy is achieved by performing one single forward pass with a single set of sampled model parameters. This is similar to the Thompson sampling strategy used by Gal and Ghahramani [9]. During training, we train the neural network parameters θ by using a Gaussian negative log-likelihood loss. The combination of DQN and dropout has been shown to practically improve sample efficiency [9]. However, dropout has multiple limitations. First, the dropout uncertainty estimates *provably* do not concentrate with more observed data [36], thus potentially violating des. 3.1. Second, there is no guarantee that dropout produce meaningful uncertainty estimates for extreme input states with a finite number of samples K , thus potentially violating des. 3.3. Third, dropout might be computationally expensive for large K value since it would require many forward passes for uncertainty estimation.

Ensemble. DQN is combined with ensembles [10] in three steps: **(1)** it trains K independent models with parameters $\{\theta_k\}_{k=1}^K$, **(2)** it performs K forward passes $\mu_k, \sigma_k = f_{\theta_k}(s^{(t)}, a^{(t)})$, and **(3)** it aggregates predictions to form the mean prediction $\mu(s^{(t)}, a^{(t)}) = \frac{1}{K} \sum_{k=1}^K \mu_k$, the aleatoric uncertainty estimate $u_{\text{alea}}(s^{(t)}, a^{(t)}) = \frac{1}{K} \sum_{k=1}^K \sigma_k$, and the epistemic uncertainty estimate $u_{\text{epist}}(s_t, a_t) = \frac{1}{K} \sum_{k=1}^K (\mu_k - \mu(s^{(t)}, a^{(t)}))^2$. In this case, the aleatoric distribution is Gaussian while the epistemic distribution is implicitly represented by the parameters of the K networks $\{\theta_k\}_{k=1}^K$. Further, the sampling-epistemic strategy is achieved by performing one single forward pass with one randomly selected network. This is similar to the Thompson sampling strategy used by Osband *et al.* [37]. We train the K independent neural network parameters θ_k with a Gaussian negative log-likelihood loss. However, ensemble has multiple limitations. First, while the combination of DQN with bootstrapped ensemble and prior functions has been *empirically* shown to improve learning for complex tasks with sparse rewards [36], [37], there is no explicit theoretical or empirical evidence that their uncertainty estimates concentrate with more observed data. Second, there is no guarantee that ensembles produce meaningful uncertainty estimates for extreme input states with a finite number of samples K , thus potentially violating des. 3.3. Third, ensemble is computationally expensive for large K value since it would require many forward passes and many neural networks.

Deep Kernel Learning. DQN is combined with deep kernel learning [14] in three steps: **(1)** it predicts one latent representation of each input state i.e. $z^{(t)} = f_{\theta}(s^{(t)})$, and **(2)** one Gaussian Process per action a defined from a fixed set of K learnable inducing points $\{\phi_{a,k}\}_{k=1}^K$ and a predefined positive definite kernel $\kappa(\cdot, \cdot)$ predicts the mean $\mu(s^{(t)}, a)$ and the variance $\sigma(s^{(t)}, a)$ of a Gaussian distribution. We train the neural network parameters θ and the inducing points $\{\phi_{a,k}\}_{k=1}^K$ jointly with a variational ELBO loss similarly to [14]. In this case, the epistemic distribution is Gaussian [13], i.e. $u_{\text{epist}}(s_t, a_t) = \mathbb{H}(\mathcal{N}(\mu(s^{(t)}, a^{(t)}), \sigma(s^{(t)}, a^{(t)})))$. Indeed, we show *theoretically* that epistemic uncertainty increases far from training data (see app. A). Thus, the combination of DQN and deep kernel learning does not suffer from arbitrary uncertainty estimates for extreme input states contrary to ReLU networks [35]. However, one of the limitation of deep kernel learning is that it does not disentangle *aleatoric* and *epistemic* uncertainty. This is similar to deep kernel learning methods in SL [13]–[15].

Evidential Networks. The combination of DQN and the posterior networks [17], [19] which belong to the class of evidential networks consists in three steps: **(1)** an encoder f_{θ} predicts one latent representation of each input state i.e. $z^{(t)} = f_{\theta}(s^{(t)})$, **(2)** one normalizing flow density estimator $\mathbb{P}(\cdot | \omega_a)$ and one linear decoder g_{ψ_a} per action a predict a Normal Inverse-Gamma distribution $\mathbb{Q}(\chi(s^{(t)}, a), n(s^{(t)}, a))$ with parameters $\chi(s^{(t)}, a) = g_{\psi_a}(z^{(t)}, a)$ and $n(s^{(t)}, a) \propto \mathbb{P}(z^{(i)} | \omega_a)$, and **(3)** it computes the posterior parameters $\chi^{\text{post}}(s^{(t)}, a) = \frac{n^{\text{prior}} \chi^{\text{prior}} + n(s^{(t)}, a) \chi(s^{(t)}, a)}{n^{\text{prior}} + n(s^{(t)}, a)}$, $n^{\text{post}}(s^{(t)}, a) = n^{\text{prior}} + n(s^{(t)}, a)$

where the prior parameters are chosen to enforce high entropy for the prior distribution e.g. $\chi^{\text{prior}} = (0, 100)^T$, $n^{\text{prior}} = 1$ [19]. In this case, the epistemic distribution is a Normal Inverse-Gamma distribution and the aleatoric distribution is a Normal distribution [19]. We train the neural network parameters θ and ψ and the normalizing flow parameters ω jointly with the MSE loss. The entropy of the conjugate prior distribution represents the epistemic uncertainty, i.e. $u_{\text{epist}}(s^{(t)}, a^{(t)}) = \mathbb{H}(\mathcal{N}\Gamma^{-1}(\chi(s^{(t)}, a^{(t)}), n(s^{(t)}, a^{(t)})))$. The entropy of the likelihood distribution represents the aleatoric uncertainty, i.e. $u_{\text{alea}}(s^{(t)}, a^{(t)}) = \mathbb{H}(\mathcal{N}(\mu(s^{(t)}, a^{(t)}), \sigma(s^{(t)}, a^{(t)})))$. Indeed, it has been showed *theoretically* that epistemic uncertainty increases far from training data (see app. A) [19]. Thus, the combination of DQN and posterior networks does not suffer from arbitrary uncertainty estimates for extreme input states contrary to ReLU networks [35].

5 Evaluation of Uncertainty Quantification in RL

In this section, we provide an extensive evaluation of uncertainty estimation for model-free RL. It compares four uncertainty estimation methods for model-free RL in three environments. First, we evaluate the uncertainty predictions at *training time* to assess the uncertainty concentration (des. 3.1) and the sample efficiency of the uncertainty-guided exploration-exploitation strategy (des. 3.2). Second, we evaluate the uncertainty estimates at *testing time* to assess the OOD detection performances (des. 3.3) and the generalization performances of the uncertainty-guided decision strategy (des. 3.4). In particular, we evaluate the trade-off between the generalization performance and the detection performance in new perturbed test environments.

Models. We consider the four uncertainty models MC Dropout (**DropOut**), **Ensemble**, deep kernel learning (**DKL**) and the evidential model based on Posterior Networks (**PostNet**) combined with the DQN RL policy (see sec. 4). In this work, we focus on DQN [22] since it is a widely used model-free RL agent which does not provide any uncertainty estimates by default. All models use the same encoder architecture and DQN hyper-parameters. We performed a grid search over all hyper-parameters. We compute the mean and standard error of the mean over 5 seeds. Further details are given in app. B.

Environments. We used three training environments **CartPole** [38], **Acrobot** [39], [40] and **LunarLander** [41], [42] from the Open AI gym environments [43]. Packer *et al.* [44] also used similar environment to assess generalization in RL. We focus on these environments since they turn out to be already challenging settings for the uncertainty methods and the sampling strategies. Some methods and strategies are indeed already unable to achieve high performance for sample efficiency, generalization, and OOD detection. We provide further details on the environments in app. C. OOD environments: The states, actions, and transition dynamics of the OOD environments should not be relevant to the original training environment task, thus being a reasonable failure mode. To this end, the input state is composed of Gaussian noise at every time step independently of the previous actions. Perturbed environments: These environments are perturbed versions of the original training environment with different perturbation strengths. We separately perturb the *state* space, the *action* space, and the *transition* dynamics with different strengths of Gaussian or uniform noises. These perturbations follow the MDP structure of the environment as proposed by the formal framework for domain shifts presented in [30]. We did not consider perturbation on the initial state only, which would be a weaker version of the state perturbations, and perturbations on the reward function which would not affect the model at testing time. Further details are given in app. C.

Training Time. First, we compare the sample efficiency and the uncertainty predictions of the four uncertainty methods using the epsilon-greedy exploration-exploitation strategy at training time. We normalize the epistemic uncertainty in $[0, 1]$ with min-max normalization to compute the relative epistemic uncertainty. It allows us to easily compare the trend of the epistemic uncertainty of all models. We show the key results for Cartpole in fig. 2 and the detailed results for CartPole, Acrobot and LunarLander in fig. 6, 7, 8 in app. F. We observe that all methods achieve similar sample efficiency. Ensemble with epsilon-greedy strategy struggles to maintain high reward on CartPole. This can be intuitively explained by the under-exploration of the epsilon-greedy strategy as also observed by Gal and Ghahramani [9] and Osband *et al.* [36]. Further, we observe that only the epistemic uncertainty estimates of the combination of DQN and PostNet decreases during training. Thus, PostNet *empirically* validates des. 3.1. In contrast, the epistemic uncertainty estimates of other methods increase or do not converge. This corroborates with the findings of [36] which *theoretically* shows that the uncertainty estimates of dropout and ensemble might not converge even on simple tasks.

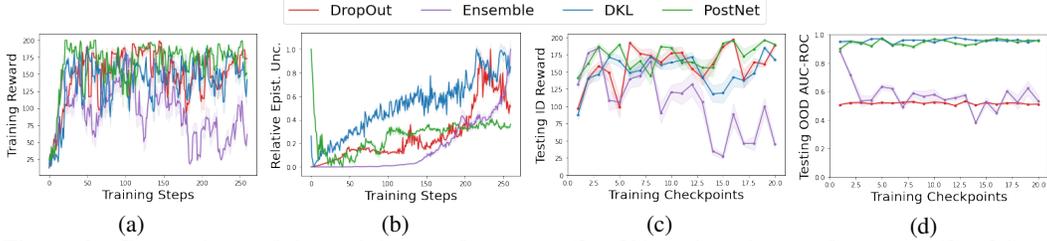


Figure 2: Comparison of the training performance (2a, 2b) and testing performance (2c, 2d) of the four uncertainty methods using epsilon-greedy strategies on CartPole. Ideally, an uncertainty aware-model in RL should achieve high reward with few training samples at training and testing time, a decreasing epistemic uncertainty at training time and high OOD detection scores at testing time.

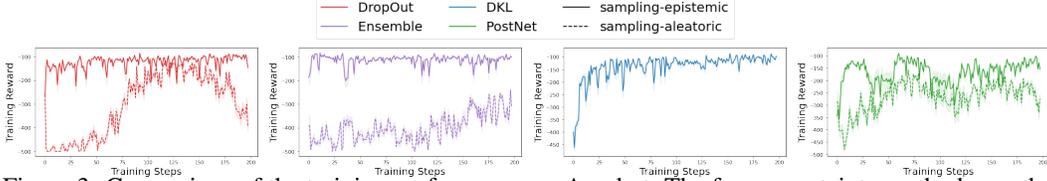


Figure 3: Comparison of the training performance on Acrobot. The four uncertainty methods use the sampling-aleatoric or the sampling-epistemic at training time. Ideally, an uncertainty aware-model should high reward with few samples.

Second, we compare the sample and episode efficiency of the *sampling-epistemic* and the *sampling-aleatoric* strategies for each model during training. We show the results for Acrobot in fig. 3, and additional results for CartPole and LunarLander in fig. 9, 11 in app. F. We observe that all models achieve high rewards by using the sampling-epistemic strategy. In particular, we observed that Ensemble with epistemic sampling achieves more stable rewards than with epsilon-greedy which aligns with observations in [37]. Further, we observe that all models instantiating both aleatoric and epistemic uncertainty achieve significantly better sample efficiency with sampling-epistemic than sampling-aleatoric. Contrary to the sampling-epistemic strategy, the sampling-aleatoric strategy intuitively fails at visiting new under-explored states/actions, thus achieving low and unstable rewards. Hence, Drop-Out, Ensemble and PostNet *empirically* validate des. 3.2. Thus, disentangling aleatoric and epistemic uncertainty can speed learning in a training environment. Further, the sampling-epistemic strategy requires fewer finished episodes on CartPole and LunarLander. This represents a more reliable training for these two environments since each finished episode translates into a failure and a restart of the systems (see app. F).

Testing Time. In this section, we save 20 models during training to evaluate their performance at testing time. The testing performance can be viewed as the model performance after deployment. First, we evaluate the testing in-distribution (ID) reward in the training environment and the out-of-distribution (OOD) detection performance against the OOD environment composed of fully noisy states. All the methods used the same epsilon-greedy strategy at training time and the action lead to the highest predicted expected return at testing time. The OOD detection performance is measured by comparing the predicted epistemic uncertainty of the states/actions of 10 episodes with the area under the receiver operating characteristic curve (AUC-ROC). We show the results for CartPole in fig. 2, and additional results for Acrobot and LunarLander in fig. 13 and fig. 14 in app. F. We observe that DKL and PostNet achieve very high OOD detection scores compared to DropOut and Ensemble. These *empirical* results align with the *theoretical* results stating that DKL and PostNet should assign high uncertainty to states very different from states observed during training. Thus, DKL and PostNet validate des. 3.3. In particular, DKL and PostNet can reliably equip DQN with epistemic uncertainty estimates which can be used to flag anomalous OOD states. In contrast, DropOut and Ensemble achieve poor OOD detection scores. This aligns with Charpentier *et al.* [19] and Osband *et al.* [36] which shows on multiple experiments that the uncertainty estimates assigned to OOD inputs by DropOut and Ensemble are not significantly smaller than the uncertainty estimates assigned to inputs close to training data.

Second, we compare the testing in-distribution (ID) reward and the out-of-distribution (OOD) detection performance when models use the *sampling-aleatoric* and *sampling-epistemic* strategies at *both* training and testing time. We show the results for the testing reward and the OOD detection scores on the LunarLander in fig. 4, and additional results on the CartPole and the Acrobot environments in fig. 15 and fig. 16 in app. F. We observe that the sampling-epistemic strategy achieves significantly better rewards than the sampling-aleatoric for almost any checkpointed models during

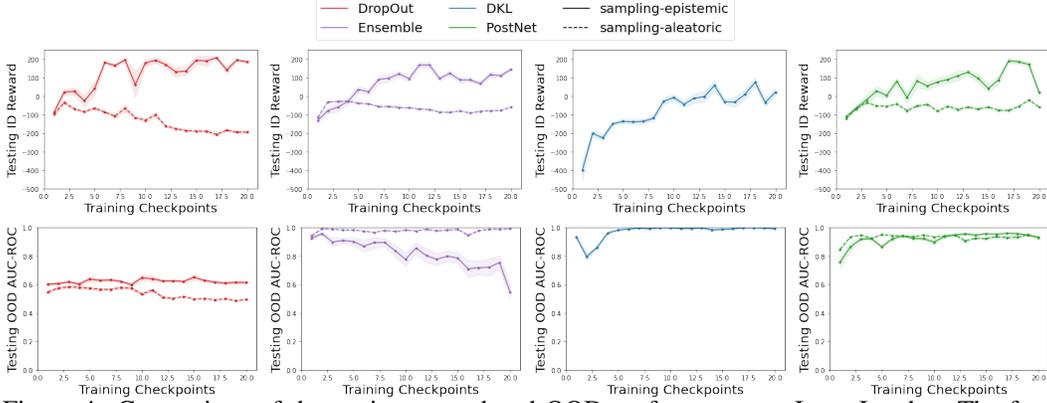


Figure 4: Comparison of the testing reward and OOD performance on LunarLander. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-ROC detection score.

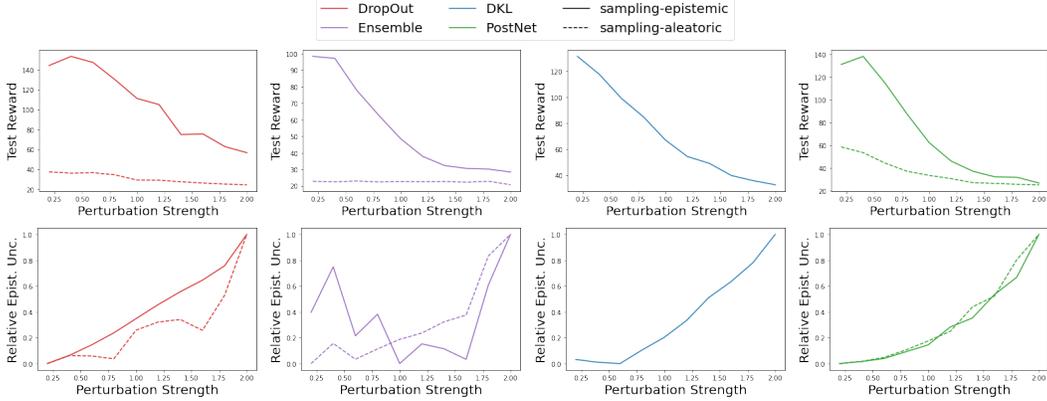


Figure 5: Comparison of the testing performance and the epistemic uncertainty predictions on CartPole with perturbed states. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

training. Thus, all models *empirically* satisfy des. 3.3. These empirical results underline the need to disentangle both aleatoric and epistemic uncertainty for high reward performance at testing time.

Third, we compare the *sampling-epistemic* and the *sampling-aleatoric* strategies for each model at testing time. All models use the same epsilon-greedy strategy at training time. We show the key results for the testing reward and the testing epistemic uncertainty on Cartpole with perturbed states in fig. 5, and detailed results with other perturbations and environments in app. F. We observe that stronger state and action perturbations deteriorate the reward performance of all models. This is reasonable since the input state or the output actions become more different from the training environment with stronger perturbations. Further, while the models were trained using the same epsilon-greedy strategy, we observe that the sampling-epistemic strategy generalizes significantly better to all types of perturbed environments than the sampling-aleatoric strategy. In particular, all models achieve high rewards with epistemic sampling on environments with perturbed transitions. Intuitively, sampling-aleatoric select actions with more inherent risk, while the sampling-epistemic select actions accounting for the knowledge accumulated by the agent in the training environment. The generalization capacity of the sampling-epistemic strategy aligns with [2] which recast the problem of generalization in RL as solving an epistemic POMDP. Thus, differentiating between aleatoric and epistemic uncertainty can improve generalization. Finally, we observed that DKL and PostNet consistently assign higher epistemic uncertainty to environments with perturbed states which aligns with their theoretical guarantees on extreme input states. The most challenging perturbations are perturbed actions since none of the models provide guarantees for this perturbation type. Overall, DKL and PostNet reliably assign higher epistemic uncertainty to most of the perturbation types. Therefore, DKL and PostNet performs a good trade-off between generalization and detection of new perturbed environments.

6 Related Work

In this section, we cover the related work for aleatoric and epistemic uncertainty estimation for RL. To this end, we review implicit *desiderata* for aleatoric and epistemic uncertainty, RL *methods* using uncertainty estimates and the existing *evaluation* methodology to validate the quality of the uncertainty estimates in practice. We refer the reader to the survey [6] for an exhaustive overview on uncertainty estimation. *Desiderata:* The notion of *risk* is well studied in RL and closely connected to the notion of uncertainty. Risk-sensitive RL usually aims at reducing the number of failures at training time for safer RL [45]–[48]. In particular, [49]–[51] discuss the trade-off between risk-sensitivity and sample efficiency at training time. Further, [2] aim at improving generalization at testing time within the Bayesian RL framework. In SL, the predicted uncertainty is expected to increase further from training data [19], [52]–[54]. None of these previous works give *explicit* desiderata for both *aleatoric* and *epistemic uncertainty* in RL at both *training* and *testing time*. *Models:* The related work for uncertainty-aware model in RL is rich as shown in surveys on distributional and Bayesian RL [26], [55]. Distributional RL [25]–[27], [56] aims at learning the distribution of return which generally captures the aleatoric uncertainty [57]. Bayesian RL methods includes sampling-based methods models such as on dropout [9], [58] and ensembles [2], [36], [37], [59], [60]. These methods are often combined with bootstrapping during training. In particular, [49] proposed to decompose aleatoric and epistemic uncertainty to the cost of multiple trained networks and [61] decompose aleatoric and epistemic uncertainty with latent variables for model-based RL. Bayesian RL also includes Gaussian processes [62], [63] and more specifically the deep kernel learning method [64] which requires storing uncertainty estimates in the experience replay buffer during training. Unlike RL, SL includes many uncertainty methods using deep kernel learning [13]–[15] and evidential network [8], [17]–[20], [54], [65]. In contrast, we look at *both* sampled-based and sampled-free uncertainty methods for *aleatoric* and *epistemic* uncertainty estimation with *minimal* modification to the training procedure of the RL agent, thus ensuring easy adaptation of new uncertainty quantification techniques from SL to RL. *Evaluation:* [66], [67] proposed to evaluate uncertainty in RL by focusing on joint predictive distributions instead of marginal distributions. Many works [1], [68]–[70] used sample efficiency as evaluation method. Further, previous works proposed generalization benchmarks for RL [44], [71]–[73]. Finally, [74], [75] have recently proposed benchmarks for OOD detection relevant to RL. In contrast, we propose a simple evaluation method which *jointly* look at *multiple* tasks relevant to real-world applications of uncertainty in RL. It covers epistemic uncertainty tracking and sample efficiency at training time, and generalization and OOD detection at testing time. In particular, we evaluate the trade-off between OOD generalization [76] and OOD detection [3].

7 Limitations and Broader Impact

Desiderata: Our desiderata, similar to [54], [77], [78], are designed to be application and model agnostic. In practice, the desiderata should be instantiated with formal definitions and could be customized depending on the application. *Models:* To validate the key contributions, similar to [25], [27], [37], we restrict our experiments to DQNs. However, the four uncertainty methods essentially modify the encoder architecture, it is possible to adapt them to other model-free RL methods such as PPO [33] and A2C [34]. *Evaluation:* Our approach focuses on a simple and task-diverse evaluation methodology for uncertainty estimation. Contrary to [73], we do not focus on scaling RL methods to more complex tasks in this paper. *Broader impact:* Our framework discusses the benefit of using uncertainty estimation to create robust and safe RL methods which corroborate with the Assessment List for Trustworthy AI [79]. Although, there is always a risk that this framework does not fully capture the real-world complexity, thus encouraging practitioners to proactively validate their models in the real-world.

8 Conclusion

We introduce a new framework to characterize aleatoric and epistemic uncertainty estimation in RL. It includes four explicit desiderata, four RL models inspired from SL and a practical evaluation methodology. The desiderata characterize the behavior of uncertainty estimates at both training and testing time. The models combine DQN with sampling-based and sampling-free uncertainty methods in SL without modifications of the RL agents training. We give theoretical and empirical evidence that these methods can fulfil the uncertainty desiderata. The evaluation method assesses the quality of uncertainty estimates on sample efficiency, generalization and OOD detection tasks.

Acknowledgments

The authors would like to thank Daniel Zügner for the helpful discussion and comments. The authors of this work take the full responsibilities for its content.

References

- [1] Z. Wang, V. Bapst, N. Heess, V. Mnih, R. Munos, K. Kavukcuoglu, and N. de Freitas, “Sample efficient actor-critic with experience replay,” *International Conference on Learning Representations*, 2016.
- [2] D. Ghosh, J. Rahme, A. Kumar, A. Zhang, R. P. Adams, and S. Levine, “Why generalization in RL is difficult: Epistemic POMDPs and implicit partial observability,” in *Advances in Neural Information Processing Systems*, 2021.
- [3] J. Yang, K. Zhou, Y. Li, and Z. Liu, “Generalized out-of-distribution detection: A survey,” *arXiv preprint arXiv:2110.11334*, 2021.
- [4] J. Nitsch, M. Itkina, R. Senanayake, J. Nieto, M. Schmidt, R. Siegwart, M. J. Kochenderfer, and C. Cadena, “Out-of-distribution detection for automotive perception,” in *24th IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2021.
- [5] Y. Gal, “Uncertainty in deep learning,” Ph.D. dissertation, University of Cambridge, 2016.
- [6] M. Abdar, F. Pourpanah, S. Hussain, D. Rezazadegan, L. Liu, M. Ghavamzadeh, P. Fieguth, A. Khosravi, U. R. Acharya, V. Makarenkov, *et al.*, “A review of uncertainty quantification in deep learning: Techniques, applications and challenges,” *Information Fusion*, vol. 76, pp. 243–297, 2021.
- [7] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. Dillon, B. Lakshminarayanan, and J. Snoek, “Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift,” in *Advances in Neural Information Processing Systems*, 2019.
- [8] A. Kopetzki, B. Charpentier, D. Zügner, S. Giri, and S. Günnemann, “Evaluating robustness of predictive uncertainty estimation: Are dirichlet-based models reliable?” *Computing Research Repository*, 2020.
- [9] Y. Gal and Z. Ghahramani, “Dropout as a bayesian approximation: Representing model uncertainty in deep learning,” in *International Conference on Machine Learning*, 2016.
- [10] B. Lakshminarayanan, A. Pritzel, and C. Blundell, “Simple and scalable predictive uncertainty estimation using deep ensembles,” in *Advances in Neural Information Processing Systems*, 2017.
- [11] F. Wenzel, J. Snoek, D. Tran, and R. Jenatton, “Hyperparameter ensembles for robustness and uncertainty quantification,” *Advances in Neural Information Processing Systems*, 2020.
- [12] Y. Wen, D. Tran, and J. Ba, “Batchensemble: An alternative approach to efficient ensemble and lifelong learning,” in *International Conference on Learning Representations*, 2020.
- [13] W. J. Maddox, P. Izmailov, T. Garipov, D. P. Vetrov, and A. G. Wilson, “A simple baseline for bayesian uncertainty in deep learning,” in *Advances in Neural Information Processing Systems*, 2019.
- [14] J. van Amersfoort, L. Smith, A. Jesson, O. Key, and Y. Gal, “On feature collapse and deep kernel learning for single forward pass uncertainty,” *arXiv preprint arXiv:2102.11409*, 2021.
- [15] J. van Amersfoort, L. Smith, Y. W. Teh, and Y. Gal, “Uncertainty estimation using a single deep deterministic neural network,” in *International Conference on Machine Learning*, 2020.
- [16] M. Biloš, B. Charpentier, and S. Günnemann, “Uncertainty on asynchronous time event prediction,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [17] B. Charpentier, D. Zügner, and S. Günnemann, “Posterior network: Uncertainty estimation without ood samples via density-based pseudo-counts,” in *Advances in Neural Information Processing Systems*, 2020.
- [18] A. Malinin and M. Gales, “Predictive uncertainty estimation via prior networks,” in *Advances in Neural Information Processing Systems*, 2018.
- [19] B. Charpentier, O. Borchert, D. Zügner, S. Geisler, and S. Günnemann, “Natural Posterior Network: Deep Bayesian Predictive Uncertainty for Exponential Family Distributions,” in *International Conference on Learning Representations*, 2022.

- [20] A. Amini, W. Schwarting, A. Soleimany, and D. Rus, “Deep evidential regression,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 14 927–14 937, 2020.
- [21] A. Malinin, N. Band, Y. Gal, M. Gales, A. Ganshin, G. Chesnokov, A. Noskov, A. Ploskonosov, L. Prokhorenkova, I. Provilkov, V. Raina, V. Raina, D. Roginskiy, M. Shmatova, P. Tigas, and B. Yangel, “Shifts: A dataset of real distributional shift across multiple large-scale tasks,” in *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [22] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” in *Advances in Neural Information Processing Systems*, 2013.
- [23] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, 2015.
- [24] W. R. Thompson, “On the likelihood that one unknown probability exceeds another in view of the evidence of two samples,” *Biometrika*, vol. 25, no. 3-4, pp. 285–294, 1933.
- [25] M. G. Bellemare, W. Dabney, and R. Munos, “A distributional perspective on reinforcement learning,” in *International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 70, 2017.
- [26] M. G. Bellemare, W. Dabney, and M. Rowland, *Distributional Reinforcement Learning*. MIT Press, 2022.
- [27] W. Dabney, G. Ostrovski, D. Silver, and R. Munos, “Implicit quantile networks for distributional reinforcement learning,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 80, 2018.
- [28] S. Agrawal and N. Goyal, “Analysis of thompson sampling for the multi-armed bandit problem,” in *Conference on learning theory*, 2012.
- [29] D. Russo and B. Van Roy, “An information-theoretic analysis of thompson sampling,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2442–2471, 2016.
- [30] T. Haider, F. S. Roza, D. Eilers, K. Roscher, and S. Günemann, “Domain shifts in reinforcement learning: Identifying disturbances in environments,” in *Workshop on AISafety at the International Joint Conference on Artificial Intelligence*, 2021.
- [31] D. Wolpert and W. Macready, “No free lunch theorems for optimization,” *IEEE Transactions on Evolutionary Computation*, 1997.
- [32] R. Bellman, “Dynamic programming,” *Science*, vol. 153, no. 3731, pp. 34–37, 1966.
- [33] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [34] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *International Conference on Machine Learning*, 2016.
- [35] M. Hein, M. Andriushchenko, and J. Bitterwolf, “Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [36] I. Osband, J. Aslanides, and A. Cassirer, “Randomized prior functions for deep reinforcement learning,” in *International Conference on Neural Information Processing Systems*, 2018.
- [37] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy, “Deep exploration via bootstrapped dqn,” in *Advances in Neural Information Processing Systems*, 2016.
- [38] A. G. Barto, R. S. Sutton, and C. W. Anderson, “Neuronlike adaptive elements that can solve difficult learning control problems,” *IEEE transactions on Systems, Man, & Cybernetics*, 1983.
- [39] R. S. Sutton, “Generalization in reinforcement learning: Successful examples using sparse coarse coding,” in *Advances in Neural Information Processing Systems*, vol. 8, 1995.
- [40] A. Geramifard, C. Dann, R. H. Klein, W. Dabney, and J. P. How, “Rlpy: A value-function-based reinforcement learning framework for education and research,” *Journal of Machine Learning Research*, 2015.
- [41] T. Brady and S. Paschall, “The challenge of safe lunar landing,” in *2010 IEEE Aerospace Conference*, 2010.

- [42] —, “The challenge of safe lunar landing,” in *2010 IEEE Aerospace Conference*, 2010.
- [43] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” *arXiv preprint arXiv:1606.01540*, 2016.
- [44] C. Packer, K. Gao, J. Kos, P. Krähenbühl, V. Koltun, and D. Song, *Assessing generalization in deep reinforcement learning*, 2018.
- [45] Y. Fei, Z. Yang, Y. Chen, Z. Wang, and Q. Xie, “Risk-sensitive reinforcement learning: Near-optimal risk-sample tradeoff in regret,” in *Advances in Neural Information Processing Systems*, 2020.
- [46] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone, “Risk-constrained reinforcement learning with percentile risk criteria,” *Journal of Machine Learning Research*, vol. 18, no. 167, pp. 1–51, 2018.
- [47] R. A. Howard and J. E. Matheson, “Risk-sensitive markov decision processes,” *Management Science*, vol. 18, no. 7, pp. 356–369, 1972.
- [48] J. García, Fern, and o Fernández, “A comprehensive survey on safe reinforcement learning,” *Journal of Machine Learning Research*, vol. 16, no. 42, pp. 1437–1480, 2015.
- [49] W. R. Clements, B. Van Delft, B.-M. Robaglia, R. B. Slaoui, and S. Toth, “Estimating risk and uncertainty in deep reinforcement learning,” *arXiv preprint arXiv:1905.09638*, 2019.
- [50] Y. Fei, Z. Yang, Y. Chen, Z. Wang, and Q. Xie, “Risk-sensitive reinforcement learning: Near-optimal risk-sample tradeoff in regret,” in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS’20, 2020.
- [51] H. Eriksson and C. Dimitrakakis, “Epistemic risk-sensitive reinforcement learning,” *ArXiv*, vol. abs/1906.06273, 2020.
- [52] A. Meinke and M. Hein, “Towards neural networks that provably know when they don’t know,” in *International Conference on Learning Representations*, 2020.
- [53] A. Kristiadi, M. Hein, and P. Hennig, “Being bayesian, even just a bit, fixes overconfidence in ReLU networks,” in *International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 119, 2020.
- [54] M. Stadler, B. Charpentier, S. Geisler, D. Zügner, and S. Günnemann, “Graph posterior network: Bayesian predictive uncertainty for node classification,” in *Advances in Neural Information Processing Systems*, 2021.
- [55] 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.
- [56] T. Morimura, M. Sugiyama, H. Kashima, H. Hachiya, and T. Tanaka, “Nonparametric return distribution approximation for reinforcement learning,” in *International Conference on Machine Learning*, 2010.
- [57] N. Nikolov, J. Kirschner, F. Berkenkamp, and A. Krause, “Information-directed exploration for deep reinforcement learning,” in *International Conference on Learning Representations*, 2019.
- [58] G. Kahn, A. Villafior, V. Pong, P. Abbeel, and S. Levine, *Uncertainty-aware reinforcement learning for collision avoidance*, 2017.
- [59] B. Lütjens, M. Everett, and J. P. How, “Safe reinforcement learning with model uncertainty estimates,” *CoRR*, vol. abs/1810.08700, 2018. arXiv: 1810.08700.
- [60] A. Tschantz, B. Millidge, A. K. Seth, and C. L. Buckley, *Reinforcement learning through active inference*, 2020. arXiv: 2002.12636.
- [61] S. Depeweg, J.-M. Hernandez-Lobato, F. Doshi-Velez, and S. Udluft, “Decomposition of uncertainty in Bayesian deep learning for efficient and risk-sensitive learning,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 80, 2018.
- [62] M. Kuss and C. Rasmussen, “Gaussian processes in reinforcement learning,” in *Advances in Neural Information Processing Systems*, vol. 16, 2003.
- [63] Y. Engel, S. Mannor, and R. Meir, “Reinforcement learning with gaussian processes,” in *International Conference on Machine Learning*, ser. ICML ’05, 2005.
- [64] J. Xuan, J. Lu, Z. Yan, and G. Zhang, “Bayesian deep reinforcement learning via deep kernel learning,” *International Journal of Computational Intelligence Systems*, vol. 12, pp. 164–171, 2018.

- [65] A. Malinin, S. Chervontsev, I. Provilkov, and M. Gales, “Regression prior networks,” *arXiv preprint arXiv:2006.11590*, 2020.
- [66] I. Osband, Y. Doron, M. Hessel, J. Aslanides, E. Sezener, A. Saraiva, K. McKinney, T. Lattimore, C. Szepesvari, S. Singh, B. V. Roy, R. Sutton, D. Silver, and H. V. Hasselt, “Behaviour suite for reinforcement learning,” in *International Conference on Learning Representations*, 2020.
- [67] I. Osband, Z. Wen, S. M. Asghari, V. Dwaracherla, B. Hao, M. Ibrahimi, D. Lawson, X. Lu, B. O’Donoghue, and B. Van Roy, *The neural testbed: Evaluating predictive distributions*, 2021.
- [68] J. Buckman, D. Hafner, G. Tucker, E. Brevdo, and H. Lee, “Sample-efficient reinforcement learning with stochastic ensemble value expansion,” in *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [69] S. Gu, T. Lillicrap, Z. Ghahramani, R. E. Turner, and S. Levine, “Q-prop: Sample-efficient policy gradient with an off-policy critic,” in *International Conference on Learning Representations*, 2017.
- [70] M. Botvinick, S. Ritter, J. Wang, Z. Kurth-Nelson, C. Blundell, and D. Hassabis, “Reinforcement learning, fast and slow,” *Trends in Cognitive Sciences*, vol. 23, 2019.
- [71] R. Kirk, A. Zhang, E. Grefenstette, and T. Rocktäschel, “A survey of generalisation in deep reinforcement learning,” 2021.
- [72] K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman, “Quantifying generalization in reinforcement learning,” in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 97, 2019.
- [73] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman, “Leveraging procedural generation to benchmark reinforcement learning,” in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 119, 2020.
- [74] M. H. Danesh and A. Fern, “Out-of-distribution dynamics detection: RL-relevant benchmarks and results,” *International Conference on Machine Learning workshop on Uncertainty in Deep Learning*, 2021.
- [75] A. P. Mohammed and M. Valdenegro-Toro, “Benchmark for out-of-distribution detection in deep reinforcement learning,” *NeurIPS Workshop on Bayesian Deep Learning*, 2021.
- [76] Z. Shen, J. Liu, Y. He, X. Zhang, R. Xu, H. Yu, and P. Cui, *Towards out-of-distribution generalization: A survey*, 2021.
- [77] R. Ashmore, R. Calinescu, and C. Paterson, “Assuring the machine learning lifecycle: Desiderata, methods, and challenges,” *ACM Comput. Surv.*, vol. 54, no. 5, 2021.
- [78] W. J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu, “Definitions, methods, and applications in interpretable machine learning,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 44, pp. 22 071–22 080, 2019.
- [79] “The assessment list for trustworthy artificial intelligence (altai) for self assessment,” *European Commission*, 2020.
- [80] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee, “Understanding deep neural networks with rectified linear units,” in *International Conference on Learning Representations*, 2018.
- [81] D. Duvenaud, “Automatic model construction with gaussian processes,” Ph.D. dissertation, University of Cambridge, 2014.
- [82] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- [83] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [84] A. Paszke, S. Gross, F. Massa, *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems*, 2019, pp. 8024–8035.
- [85] W. Falcon *et al.*, “Pytorch lightning,” *GitHub*. *Note: <https://github.com/PyTorchLightning/pytorch-lightning>*, vol. 3, 2019.
- [86] L. Biewald, *Experiment tracking with weights and biases*, Software available from wandb.com, 2020.

- [87] J. R. Gardner, G. Pleiss, D. Bindel, K. Q. Weinberger, and A. G. Wilson, “Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration,” in *Advances in Neural Information Processing Systems*, 2018.

Appendix

A Proofs

In this section, we show that deep kernel learning [14] and evidential models based on Posterior Networks [17], [19] are guaranteed to assign high epistemic uncertainty for inputs far from inputs observed during training under technical assumptions. In particular, the combination of DQN with deep kernel learning or evidential networks presented in sec. 4 are guaranteed to assign high epistemic uncertainty for extreme input states. We assume that the encoder should use ReLU activations, which is common in deep learning, and that the rows of the linear transformations are independent, which is realistic for trained networks with no constant output [35].

Lemma 1. [80] *Let $\{Q_l\}_l^R$ be the set of linear regions associated to the piecewise ReLU network $f_\phi(\mathbf{x})$. For any $\mathbf{x} \in \mathbb{R}^D$, there exists $\delta^* \in \mathbb{R}^+$ and $l^* \in 1, \dots, R$ such that $\delta \cdot \mathbf{x} \in Q_{l^*}$ for all $\delta > \delta^*$.*

Lemma 2. *Let a (deep) encoder f_ϕ with piecewise ReLU activations. Let $f_\phi(\mathbf{x}) = V^{(l)}\mathbf{x} + a^{(l)}$ be the piecewise affine representation of the ReLU network f_ϕ on the finite number of affine regions $Q^{(l)}$ [80]. Suppose that $V^{(l)}$ have independent rows, then for almost any \mathbf{x} we have $\|f_\phi(\delta \cdot \mathbf{x})\| \xrightarrow{\delta \rightarrow \infty} \infty$. i.e the norm of the latent representations $\mathbf{z}_\delta = f_\phi(\delta \cdot \mathbf{x})$ associated to the input $\delta \cdot \mathbf{x}$ goes to infinity.*

Proof. We prove now lem. 2. Let $\mathbf{x} \in \mathbb{R}^D$ be a non-zero input and f_ϕ be a ReLU network. Lem. 1 implies that there exists $\delta^* \in \mathbb{R}^+$ and $l \in \{1, \dots, R\}$ such that $\delta \cdot \mathbf{x} \in Q^{(l)}$ for all $\delta > \delta^*$. Thus, $\mathbf{z}_\delta = f_\phi(\delta \cdot \mathbf{x}) = \delta \cdot (V^{(l)}\mathbf{x}) + a^{(l)}$ for all $\delta > \delta^*$. Note that for $\delta \in [\delta^*, +\infty]$, \mathbf{z}_δ follows an affine half line $S_{\mathbf{x}} = \{\mathbf{z} \mid \mathbf{z} = \delta \cdot (V^{(l)}\mathbf{x}) + a^{(l)}, \delta > \delta^*\}$ in the latent space. Further, note that $V^{(l)}\mathbf{x} \neq 0$ since $\mathbf{x} \neq 0$ and $V^{(l)}$ has independent rows. Therefore, we have $\|\mathbf{z}_\delta\| \xrightarrow{\delta \rightarrow \infty} +\infty$ \square

Theorem 3. *Let a Deep Kernel Learning model parametrized with a (deep) encoder f_ϕ with piecewise ReLU activations, a set of K inducing points $\{\phi_k\}_{k=1}^K$ and a RBF, Matern or Rational Quadratic kernel $\kappa(\cdot, \cdot)$ [81], [82]. Let $f_\phi(\mathbf{x}) = V^{(l)}\mathbf{x} + a^{(l)}$ be the piecewise affine representation of the ReLU network f_ϕ on the finite number of affine regions $Q^{(l)}$ [80]. Suppose that $V^{(l)}$ have independent rows, then for almost any \mathbf{x} we have $\sigma(f_\phi(\delta \cdot \mathbf{x})) \xrightarrow{\delta \rightarrow \infty} c$ where $c = \kappa(0, 0)$.*

Proof. We prove now thm. 3. Lem. 1 says that $\|\mathbf{z}_\delta\| \xrightarrow{\delta \rightarrow \infty} +\infty$ where $\mathbf{z}_\delta = f_\phi(\delta \cdot \mathbf{x})$. It implies that $\|\mathbf{z}_\delta - \phi_k\| \xrightarrow{\delta \rightarrow \infty} \infty$ for all inducing point ϕ_k . Thus, we obtain $\kappa(\mathbf{z}_\delta, \phi_k) \xrightarrow{\delta \rightarrow \infty} 0$ where $\kappa(\cdot, \cdot)$ is the RBF, Matern or Rational Quadratic kernel [81], [82]. Since the variance of the predictive Gaussian distribution associated with the Gaussian process is $\sigma(f_\phi(\delta \cdot \mathbf{x})) = c - \boldsymbol{\kappa} \mathbf{C} \boldsymbol{\kappa}$ where $c = \kappa(f_\phi(\delta \cdot \mathbf{x}), f_\phi(\delta \cdot \mathbf{x})) = \kappa(0, 0)$, $\boldsymbol{\kappa}_k = \kappa(\mathbf{z}_\delta, \phi_k)$ and $\boldsymbol{\kappa}_{k,k'} = \kappa(\phi_k, \phi_{k'})$. This gives the final result $\sigma(f_\phi(\delta \cdot \mathbf{x})) \xrightarrow{\delta \rightarrow \infty} c$ where $c = \kappa(0, 0)$. \square

Thm. 3 implies that deep kernel learning on a latent space parametrized with a neural network is guaranteed to predict high uncertainty corresponding to the prior uncertainty far from training data. This includes the uncertainty predicted by the GP associated to each action a in the combination of DQN and deep kernel learning presented in sec. 4. The uncertainty prediction $u_{\text{epist}}(s_t, a_t) = \mathbb{H}(\mathcal{N}(\mu(\mathbf{s}^{(t)}), a^{(t)}), \sigma(\mathbf{s}^{(t)}, a^{(t)}))$ becomes high for input states $s^{(t)}$ extremely different from the training environment i.e. $\|s^{(t)}\| \rightarrow \infty$.

Theorem 4. [19] *Let a Natural Posterior Network model parametrized with a (deep) encoder f_ϕ with piecewise ReLU activations, a decoder g_ψ and the density $\mathbb{P}(\mathbf{z} \mid \boldsymbol{\omega})$. Let $f_\phi(\mathbf{x}) = V^{(l)}\mathbf{x} + a^{(l)}$ be the piecewise affine representation of the ReLU network f_ϕ on the finite number of affine regions $Q^{(l)}$ [80]. Suppose that $V^{(l)}$ have independent rows and the density function $\mathbb{P}(\mathbf{z} \mid \boldsymbol{\omega})$ has bounded derivatives, then for almost any \mathbf{x} we have $\mathbb{P}(f_\phi(\delta \cdot \mathbf{x}) \mid \boldsymbol{\omega}) \xrightarrow{\delta \rightarrow \infty} 0$. i.e the evidence becomes small far from training data.*

The proof of thm. 4 is given in [19], and relies also on lem. 2 and the fact that a smooth density estimator should converge to 0 far from training data. Intuitively, it implies that the epistemic associated to each possible action a by the combination of DQN and posterior network becomes high

for input states $s^{(t)}$ extremely different from the training environment i.e. $\|s^{(t)}\| \rightarrow \infty$. In particular, prior parameter takes over in the posterior update (i.e. $n^{\text{post}}(s^{(t)}, a) \rightarrow n^{\text{prior}}$, $\chi^{\text{post}}(s^{(t)}, a) \rightarrow \chi^{\text{prior}}$)

B Model Details

We train all models on a single GPU (NVIDIA GTX 1080 Ti or NVIDIA GTX 2080 Ti, 11 GB memory). All models use the same core architecture. They use a 2 layers MLP with 128 hidden units for the CartPole environment, a 2 layers MLP with 64 hidden units for the Acrobot environment and a 3 layers MLP with 128 hidden units for the LunarLander environment. All models are trained using 5 random seeds with the Adam optimizer [83]. For fair comparison, we use the same hyperparameters for the DQN architecture in all uncertainty models: the target network parameters are completely updated (i.e. $\tau = 1.$) every 10 training iterations. The epsilon-greedy strategy start with $\epsilon = 1.$ and decay till $\epsilon = 0.01$ after 1000 iteration steps. The discount factor is set to 0.99. Further, we use a batch size of 16, a replay size of 1000 and a maximum number of training iterations of 13000 for Cartpole, a batch size of 64, a replay size of 10000 and a maximum number of training iterations of 120000 for Acrobot, and a batch size of 128, a replay size of 10000 and a maximum number of training iterations of 300000 for LunarLander. For each type of uncertainty model, we performed a grid search for the learning rate in the range $[10^{-1}, 10^{-4}]$.

Each uncertainty method has also its own hyperparameters. We show an hyperparameter study in app. F.4 for the main hyper-parameters of each uncertainty method. For the MC dropout model, we make a grid-search over the number of samples $n \in [10, 20, 40, 80]$ and the drop probability $p \in [.1, .2, .3, .4, .5]$. In the main experiments, we use $n = 80$ and $p = .2$. For the ensemble model, we make a grid-search over the number of networks $n \in [10, 20, 40, 80]$. In the main experiments, we use $n = 80$. For the deep kernel learning model, we make a grid-search over the number of inducing points $n \in [10, 20, 40, 80]$, the latent dimension $H \in [16, 32, 64]$, the kernel type in RBF, RQ and Matern- $\frac{3}{2}$ Kernel, and an ELBO regularization factor in $\lambda \in [., 1.]$. In the main experiments, we use $n = 80$ inducing points, a latent dimension of $H = 64$, the RQ kernel, and a regularization factor of $\lambda = .1$. Further, we observed that adding a batch normalization layer right after the encoder f_{θ} was stabilizing the training similarly to Charpentier *et al.* [17]. For the evidential network model based on posterior networks, we make a grid-search over the flow depth $d \in [8, 16, 32]$, the latent dimension $H \in [8, 16, 32]$. In the main experiments, we use a radial flow with depth $d = 8$ and a latent dimension of $H = 16$. Further, we observed that adding a batch normalization layer right after the encoder f_{θ} was stabilizing the training similarly to Charpentier *et al.* [17].

We will provide the github repository with the code on the project page <https://www.cs.cit.tum.de/daml/aleatoric-epistemic-uncertainty-rl/>. To conduct the experiments, we used Pytorch [84] with BSD license, Pytorch Lightning [85] with Apache 2.0 license and Weight&Biases [86]. Further, we also use GPytorch for to implement the deep kernel model [87].

C Environment Details

We use OpenAI gym environments [43] with MIT license. We design the OOD environments such that they should not be relevant to the original training environment task, and thus being a reasonable failure mode. Further, we design a continuum of perturbed environments going from tasks very similar to the training environment to the tasks very different from the original environment. We distinguish between perturbations on the *state* space, the *action* space, and the *transition* dynamics to follow the MDP structure of the original environment. In contrast, [44], [75] mostly focus on perturbations on the environment parameters. We will provide the github repository with the code for the OOD and perturbed environment on the project page <https://www.cs.cit.tum.de/daml/aleatoric-epistemic-uncertainty-rl/>.

Cartpole [38] In this environment, the goal of the agent is to maintain a pole on a cart straight up. This environment has a discrete action space with 2 possible actions corresponding to apply the a force to the left or the right of the cart. This environment has a continuous state space with dimension 4 corresponds to. The episode ends when the pole is more than 15 degrees from vertical, or the cart moves more than 2.4 units from the center. The reward is +1 at every time step that the pole stays up. The maximum length of an episode is 200 steps. For the OOD environment, the input states are drawn from a Gaussian distribution with unit variance i.e. $s^{(t),\text{pert}} \sim \mathcal{N}(0, 1)$. For

the perturbed environments with perturbation strength ϵ , the action space is perturbed by randomly adding Gaussian noise to the scale of the force applied to the cart (i.e. $f^{\text{pert}} = (1. + x)f$ where $f \in \{-10, 10\}$ is default action force and $x \sim \mathcal{N}(0, \epsilon)$ is the perturbation), the state space is perturbed by adding Gaussian noise to the observation scale (i.e. $s^{(t),\text{pert}} = (1. + x)s^{(t)}$ where $x \sim \mathcal{N}(0, \epsilon)$ is the perturbation), and the transition dynamic is perturbed by adding a uniform noise centered around the true dynamic parameters (i.e. $\nu^{\text{pert}} = (1 + x)\nu$ where $x \sim U(-\epsilon, \epsilon)$ where ν is the original environment parameters) such as the gravity, pole length ...etc.

Acrobot [39], [40] In this environment, the agent control a robot arm with two links and its goal is to move the end of the lower link up to a given height. This environment has a discrete action space with 3 possible actions corresponding to apply a positive torque, a negative torque or nothing. This environment has a continuous state space with dimension 6 corresponding to the 4 joint angles and the 2 angular velocities. The episode ends when the lower link of the robot arm is above a given height. The reward is -1 at every time step that the pole does not reach the expected height. The maximum length of an episode is 500 steps at most. For the OOD environment, the input states are drawn from a Gaussian distribution with unit variance (i.e. $s^{(t),\text{pert}} \sim \mathcal{N}(0, 1)$). For the perturbed environments with perturbation strength ϵ , the action space is perturbed by randomly sampling actions with probability $p = \frac{\epsilon}{2}$, the state space is perturbed by adding Gaussian noise to the observation scale (i.e. $s^{(t),\text{pert}} = (1. + x)s^{(t)}$ where $x \sim \mathcal{N}(0, \epsilon)$), and the transition dynamic is perturbed by adding a uniform noise centered around the true dynamic parameters (i.e. $\nu^{\text{pert}} = (1 + x)\nu$ where $x \sim U(-\epsilon, \epsilon)$ where ν is the original environment parameters) such as the lengths of the links, the masses of the links.

LunarLander [41], [42] In this environment, the agent control a space ship and its goal is to land it on the surface of the moon. This environment has a discrete action space with 4 possible actions corresponding to apply a torque to the left, to the right, downward or nothing. This environment has a continuous state space with dimension 8 corresponding to the space ship coordinates. The reward is correlated with fast landing in the correct area without crashes. The episode ends when the spaceship is landed or crashed. For the OOD environment, the input states are drawn from a Gaussian distribution with unit variance (i.e. $s^{(t),\text{pert}} \sim \mathcal{N}(0, 1)$). For the perturbed environments with perturbation strength ϵ , the action space is perturbed by randomly sampling actions with probability $p = \frac{\epsilon}{2}$, the state space is perturbed by adding Gaussian noise to the observation scale (i.e. $s^{(t),\text{pert}} = (1. + x)s^{(t)}$ where $x \sim \mathcal{N}(0, \epsilon)$), and the transition dynamic is perturbed by adding a uniform noise centered around the true dynamic parameters (i.e. $\nu^{\text{pert}} = (1 + x)\nu$ where $x \sim U(-\epsilon, \epsilon)$ where ν is the original environment parameters) such as the lengths of the links, the masses of the links.

D Metric Details - Training Time

We track the current reward, the epistemic uncertainty and the aleatoric uncertainty at every training step. The epistemic and aleatoric uncertainty are defined by the variance or the entropy of the epistemic and the aleatoric distributions (see sec. 4). The two uncertainty types are then normalized between $[0, 1]$ with min-max normalization to compute the relative epistemic and aleatoric uncertainty on the plots. The normalization enable an easier comparison of the trend of the uncertainty estimates across methods. For all these experiments, we compute the mean and the standard error of the mean across 5 seeds for all results.

E Metric Details - Testing Time

We save 20 model checkpoints at regular interval during the whole training. We evaluate then the 20 checkpointed models at testing time. First, we compute the in-distribution (ID) reward average over 10 episodes on the original training environment. Second, we compute the OOD detection scores by comparing the epistemic uncertainty of the ID and the OOD environment over 10 episodes each with the area under the receiver operating characteristic curve (AUC-ROC) and the area under the precision-recall curve (AUC-PR). Higher scores indicate better OOD detection performances. Third, we compute the averaged reward and epistemic uncertainty on the perturbed environment over 10 episodes. For all these experiments, we compute the mean and the standard error of the mean across 5 seeds for all results. Further, we also sampled 5 random perturbations for each perturbation strength.

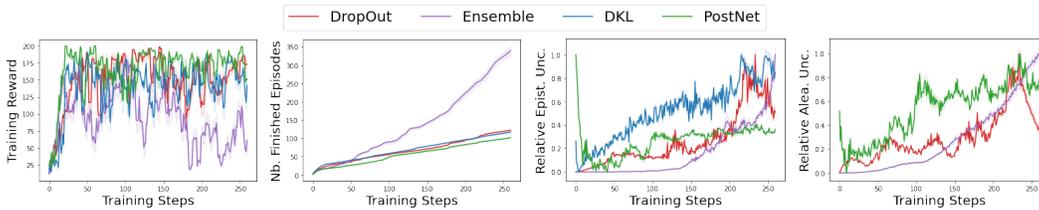


Figure 6: Comparison of the training performance of the four uncertainty methods using epsilon-greedy strategies on CartPole. Ideally, a uncertainty aware-model should achieve high reward with few samples and episodes and with a decreasing epistemic uncertainty.

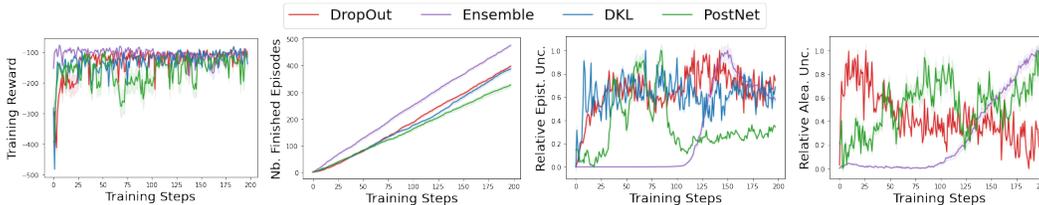


Figure 7: Comparison of the training performance of the four uncertainty methods using epsilon-greedy strategies on Acrobot. Ideally, a uncertainty aware-model should achieve high reward with few samples and episodes and with a decreasing epistemic uncertainty.

F Additional Experiments

F.1 Training Time

We show additional results on CartPole, Acrobot and LunarLander in fig. 6, fig. 7 and fig. 8 to compare the performance of the uncertainty estimates of the four uncertainty methods at training time. The epistemic uncertainty estimates of PostNet decrease during training. Thus, PostNet empirically validate des. 3.1. Further, Ensemble and PostNet require a low number of finished episodes on CartPole and LunarLander. This translates for these two environments into a safer learning with a lower number of restart of the systems.

We show additional results in fig. 9, fig. 10 and fig. 11 to compare the performance of the sampling-epistemic and the sampling-aleatoric strategies at training time. The sampling-epistemic strategy consistently achieve a better sample efficiency. Thus, Ensemble, Dropout and PostNet empirically satisfy des. 3.2. Hence, disentangling aleatoric and epistemic uncertainty can speed learning in a training environment.

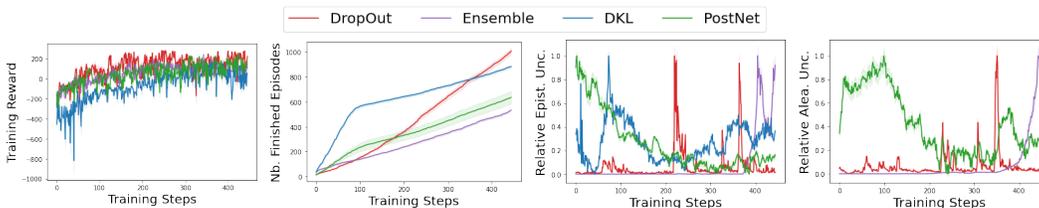


Figure 8: Comparison of the training performance of the four uncertainty methods using epsilon-greedy strategies on LunarLander. Ideally, a uncertainty aware-model should achieve high reward with few samples and episodes and with a decreasing epistemic uncertainty.

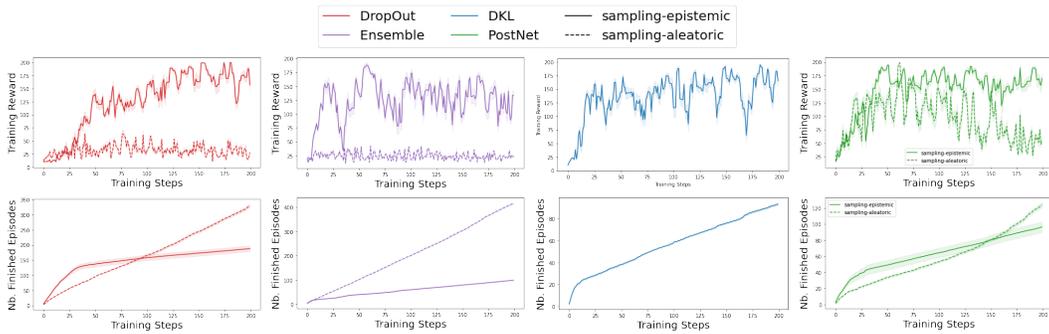


Figure 9: Comparison of the training performance on Cartpole. The four uncertainty methods use the sampling-aleatoric or the sampling-epistemic at training time. Ideally, an uncertainty aware-model should high reward with few samples.

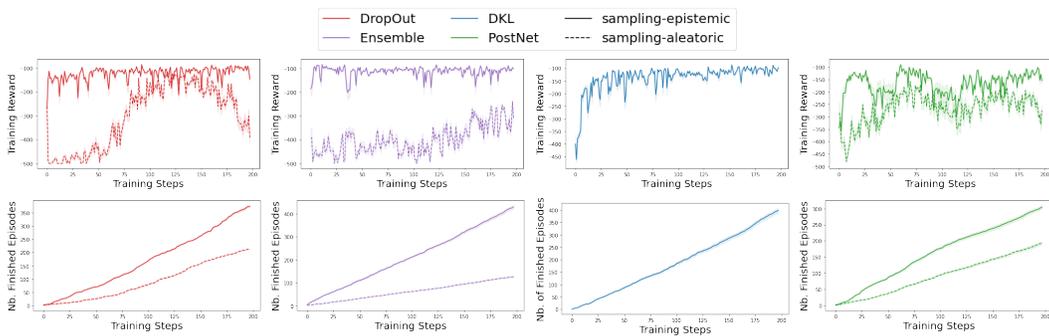


Figure 10: Comparison of the training performance on Acrobot. The four uncertainty methods use the sampling-aleatoric or the sampling-epistemic at training time. Ideally, an uncertainty aware-model should high reward with few samples.

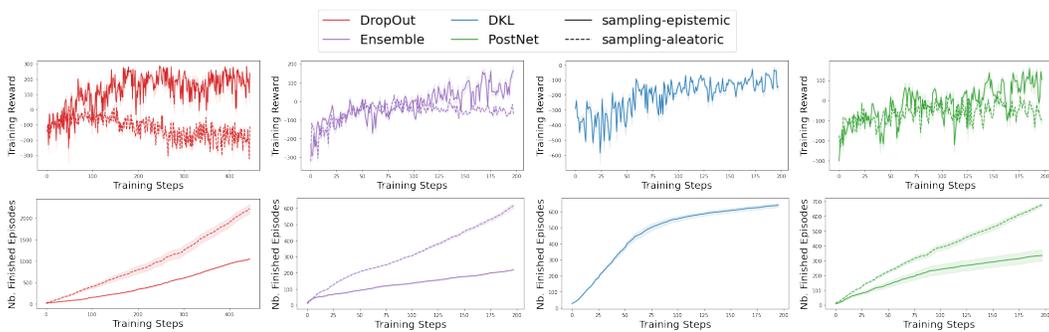


Figure 11: Comparison of the training performance on LunarLander. The four uncertainty methods use the sampling-aleatoric or the sampling-epistemic at training time. Ideally, an uncertainty aware-model should high reward with few samples.

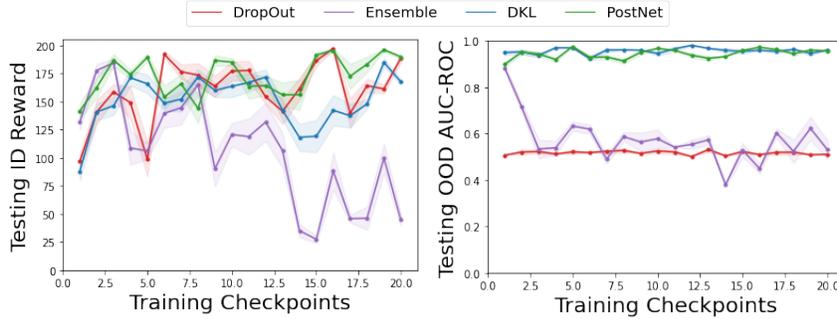


Figure 12: Comparison of the testing performance of the four uncertainty methods using epsilon-greedy strategies at training and testing time on CartPole. Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

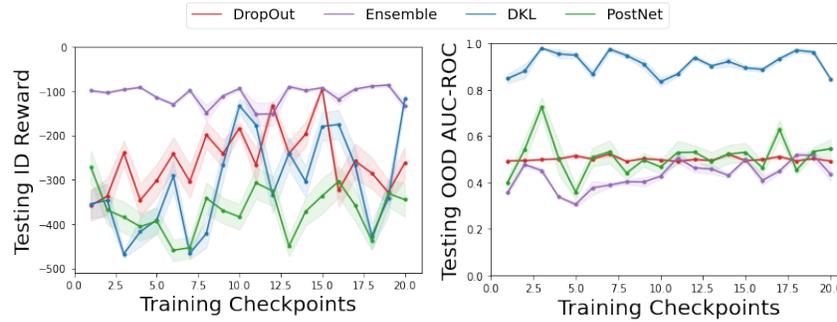


Figure 13: Comparison of the testing performance of the four uncertainty methods using epsilon-greedy strategies at training and testing time on Acrobot. Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

F.2 Testing Time

We show additional results in fig. 12, fig. 13 and fig. 14 to compare the generalization and OOD detection performance of the uncertainty estimates of the four uncertainty methods at testing time. The models use the sampling-epistemic or the sampling-aleatoric strategy at both training and testing time. Further, we show other additional results for OOD detection by using the area under the precision-recall (AUC-PR) scores instead of the area under the receiver operating characteristic curve (AUC-ROC) in fig. 25, fig. 26, fig. 27. We observe that DKL and PostNet achieve very high OOD detection scores in most settings compared to DropOut and Ensemble. These *empirical* results align with the *theoretical* results stating that DKL and PostNet should assign high uncertainty to states very different from states observed during training. Thus, DKL and PostNet validate des. 3.3. In particular, DKL and PostNet can reliably equip DQN with epistemic uncertainty estimates which can be used to flag anomalous OOD states.

We show additional results in fig. 15, fig. 16 and fig. 4 to compare the performance of the sampling-epistemic and sampling-aleatoric strategies for each uncertainty model. All models use the same epsilon-greedy strategy at training time. We observe that the sampling-epistemic strategy is consistently better than sampling-aleatoric at testing time. The higher generalization capacity of the sampling-epistemic strategy aligns with [2] which recasts the problem of generalization in RL as solving an epistemic POMDP. These empirical results underline the need to disentangle both aleatoric and epistemic uncertainty for high reward performance at testing time.

We show additional results in fig. 5, fig. 19, fig. 22, fig. 17, fig. 20, fig. 23, fig. 18, fig. 21, fig. 24 to compare the generalization and uncertainty performances of the sampling-epistemic and sampling-aleatoric strategies of each method on perturbed environments with state, action and transition dynamic perturbations. All methods achieve lower reward on environment with stronger perturbations. This is expected since a model cannot generalize to all new environments. The sampling-epistemic strategy achieves significantly better than the sampling-aleatoric strategy. The generalization capacity of the sampling-epistemic strategy aligns again with [2]. Thus, differentiating between aleatoric and

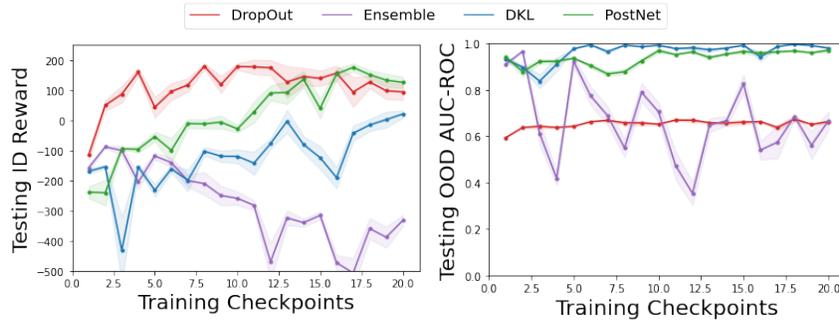


Figure 14: Comparison of the testing performance of the four uncertainty methods using epsilon-greedy strategies at training and testing time on LunarLander. Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

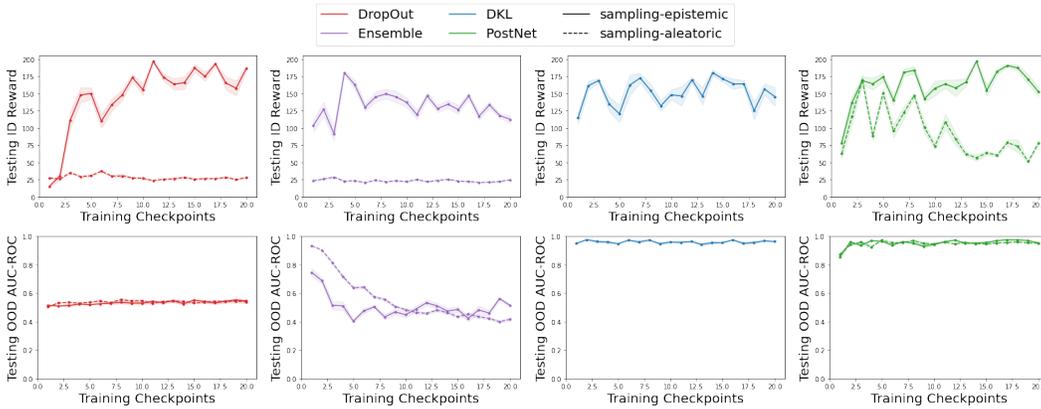


Figure 15: Comparison of the testing reward and OOD performance on CartPole. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-ROC detection score.

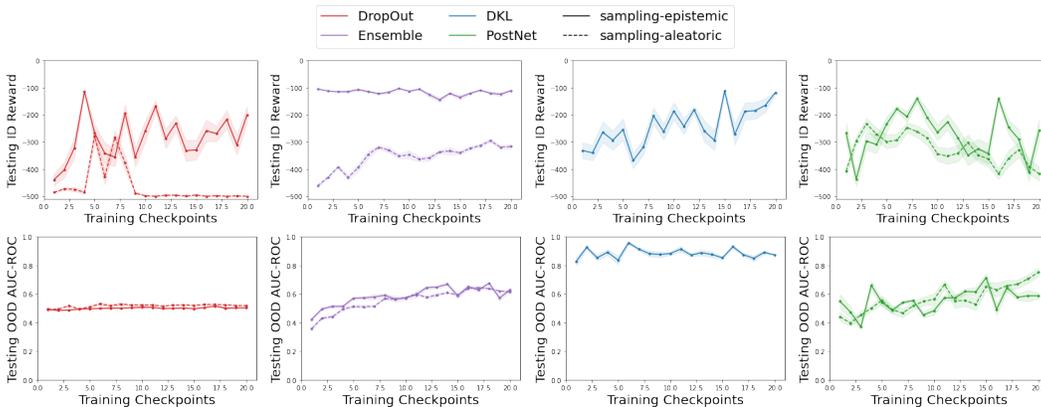


Figure 16: Comparison of the testing reward and OOD performance on Acrobot. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-ROC detection score.

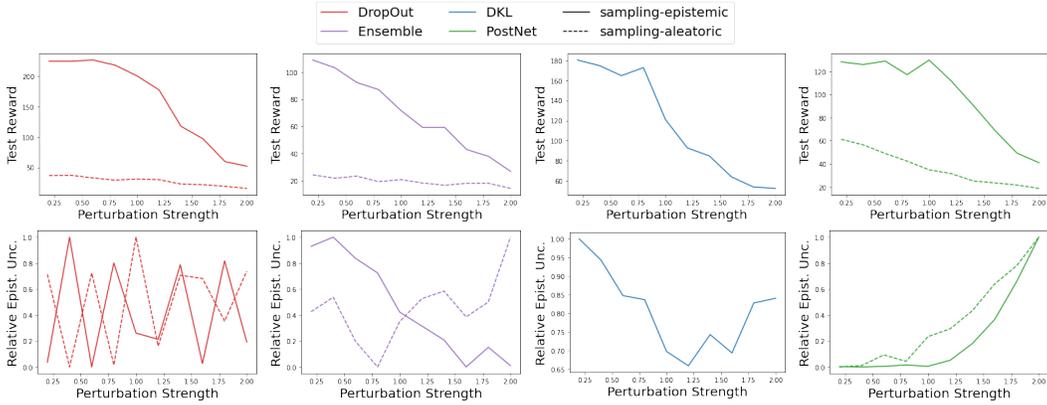


Figure 17: Comparison of the testing performance and the epistemic uncertainty predictions on CartPole with perturbed actions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

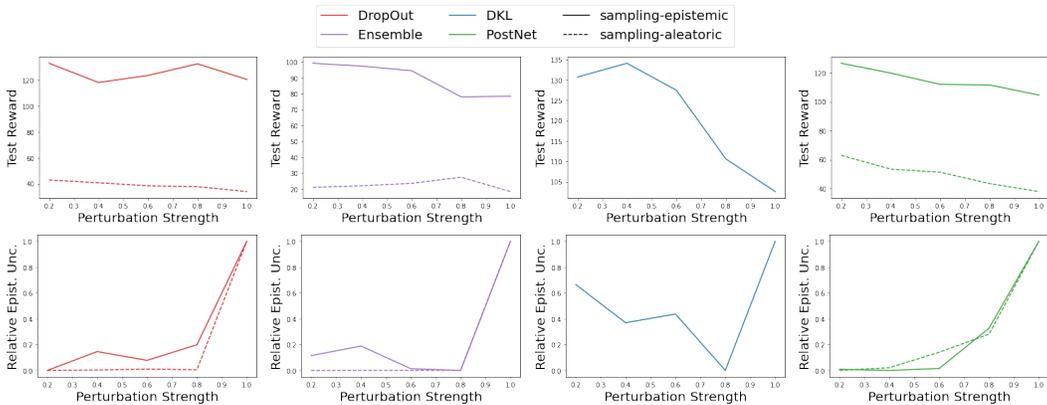


Figure 18: Comparison of the testing performance and the epistemic uncertainty predictions on CartPole with perturbed transitions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

epistemic uncertainty can improve generalization. Finally, only DKL and PostNet reliably assign higher epistemic uncertainty to most of the perturbation types. Therefore, DKL and PostNet have a good trade-off between generalization and detection of new perturbed environments.

Video: For a better visualization, we provide supplementary videos on the project page <https://www.cs.cit.tum.de/daml/aleatoric-epistemic-uncertainty-rl/>. The videos show the landing performance, the reward performance, and the relative epistemic uncertainty prediction of the PostNet model in the original LunarLander environments and two environments with perturbed states with perturbation strengths equal to 0.5 and 2.0. On the original environment, we observe that the space ship lands correctly with lower epistemic uncertainty after landing. On the perturbed environment with strength 0.5, we observe that the space ship avoids crashing but assigns higher epistemic uncertainty when moving further from the landing zone. Finally, on the perturbed environment with strength 2.0, we observe that the space ship assigns significantly higher epistemic uncertainty especially when approaching the floor before the crash.

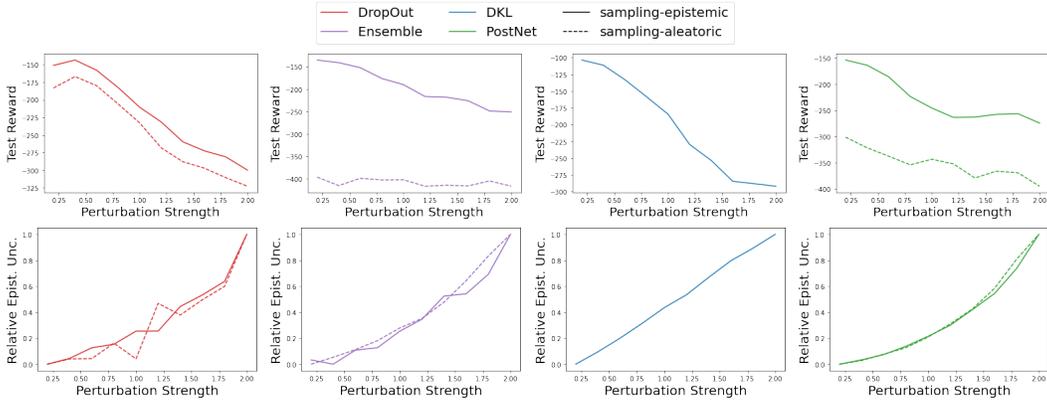


Figure 19: Comparison of the testing performance and the epistemic uncertainty predictions on Acrobot with perturbed states. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

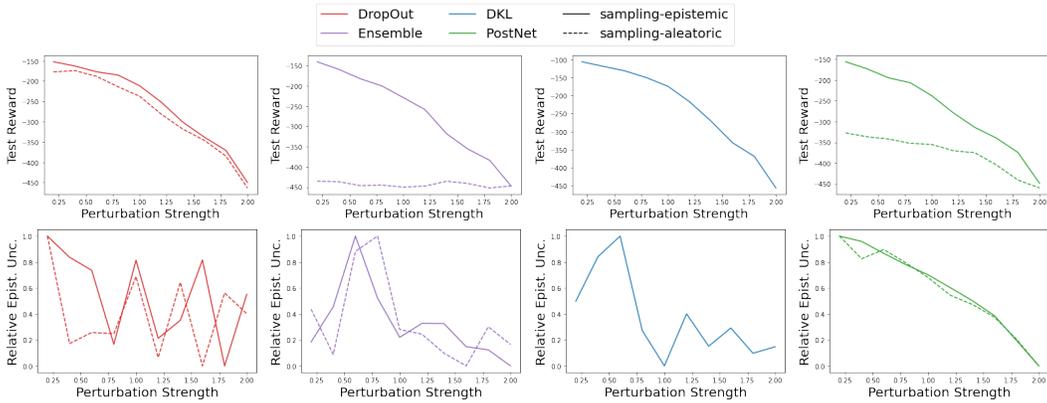


Figure 20: Comparison of the testing performance and the epistemic uncertainty predictions on Acrobot with perturbed actions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

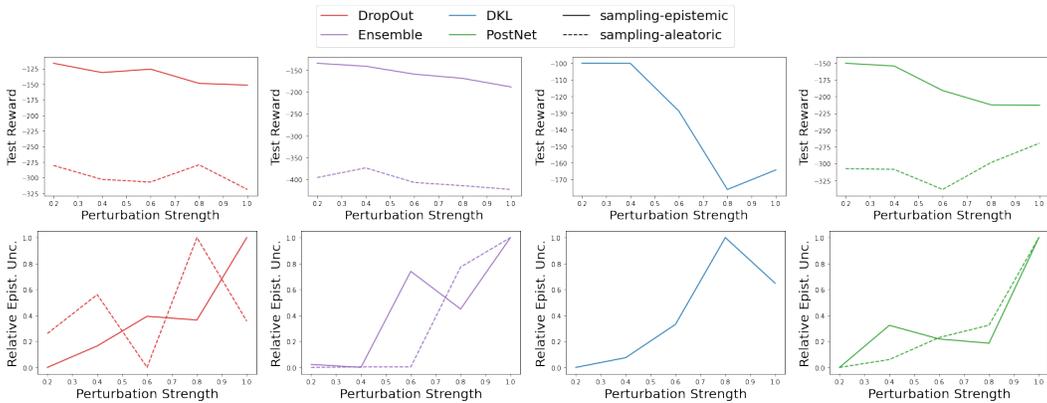


Figure 21: Comparison of the testing performance and the epistemic uncertainty predictions on Acrobot with perturbed transitions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

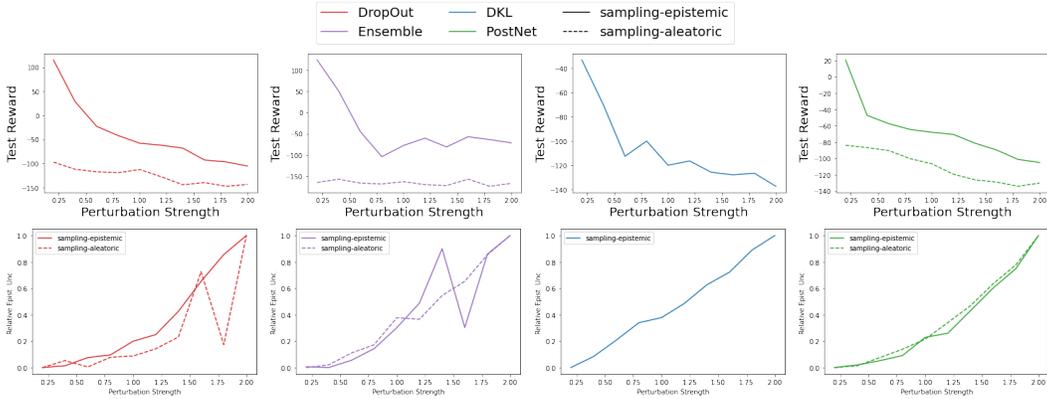


Figure 22: Comparison of the testing performance and the epistemic uncertainty predictions on LunarLander with perturbed states. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

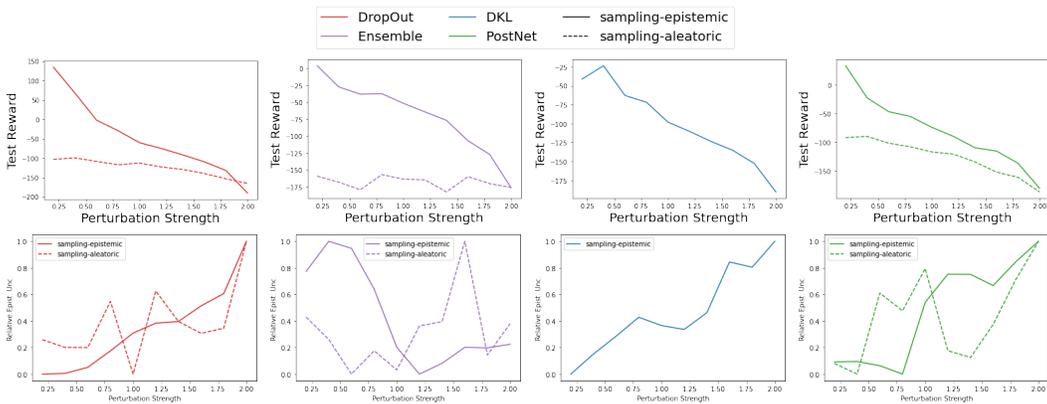


Figure 23: Comparison of the testing performance and the epistemic uncertainty predictions on LunarLander with perturbed actions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

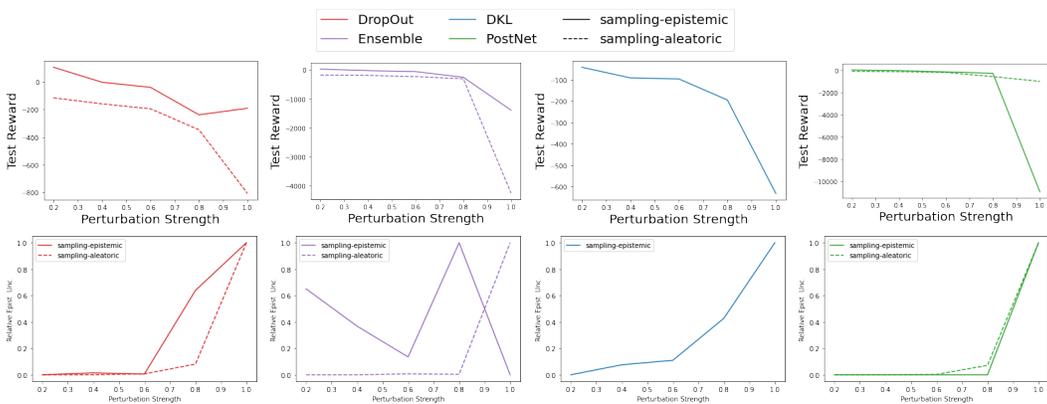


Figure 24: Comparison of the testing performance and the epistemic uncertainty predictions on LunarLander with perturbed transitions. The four uncertainty methods use the epsilon-greedy strategy at training time and the sampling-aleatoric or sampling-epistemic strategy at testing time. Ideally, an uncertainty-aware model should maintain high reward while assigning higher epistemic uncertainty on more severe perturbations.

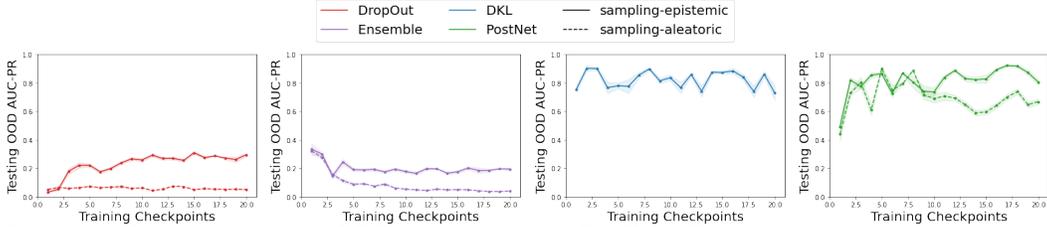


Figure 25: Comparison of the OOD performance on CartPole. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-PR detection score.

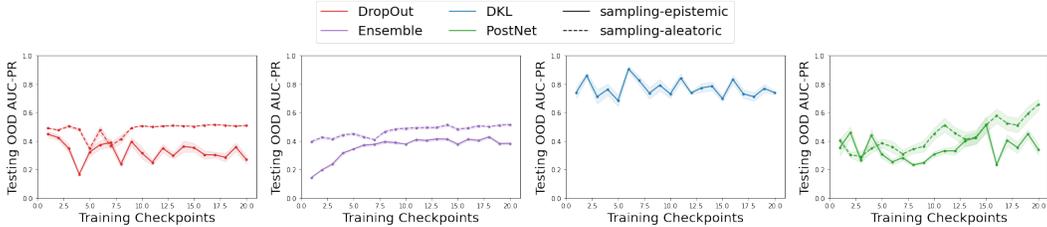


Figure 26: Comparison of the OOD performance on Acrobot. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-PR detection score.

E.3 Comparison with Vanilla DQN

We show additional results in fig. 28, fig. 29, fig. 30 to compare the sample efficiency and the generalization capacity of the uncertainty models with the vanilla DQN. The vanilla DQN is not equipped by default with uncertainty estimates. Therefore, it cannot be used for uncertainty tasks like OOD detection. For the sake of comparison, all models use the epsilon-greedy strategy. We observe that the vanilla DQN achieve significantly lower sample efficiency on CartPole. Further, it achieves less stable generalization performance on LunarLander. In contrast, the four uncertainty methods achieve higher generalization performance especially when using the sampling epistemic strategy (see fig. 15, fig. 16 and fig. 4). These results underline the benefit of predicting and disentangling the aleatoric and the epistemic uncertainty for better sample efficiency and generalization performance.

E.4 Hyperparameter Selection

In this section, we present a hyperparameter study for each uncertainty method on the CartPole environment. To this end, plot the testing reward and the OOD scores when varying the most important hyper-parameters. we at testing time. We show the hyper-parameter study for dropout when varying the number of samples n and the dropout probability p in fig. 31. We observe that a higher number of samples achieves a slightly better OOD detection score. Dropout is pretty insensitive to the dropout probability. We show the hyper-parameter study for ensemble when varying the number of networks n in fig. 32. While a higher number of networks is supposed to give higher prediction quality [10], Ensemble looks to give similar results for all number of networks. We show the hyper-parameter study for DKL when varying the number of inducing points n , the latent dimension H , the kernel type and the batch norm layer in fig. 33. The batch norm layer appears to improve the results similarly to [17].

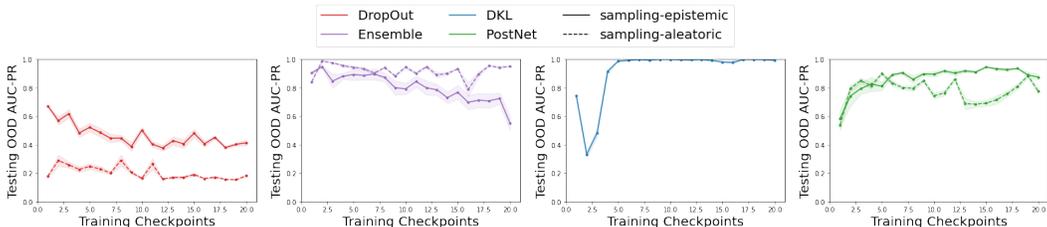


Figure 27: Comparison of the OOD performance on LunarLander. The four uncertainty methods use the sampling-aleatoric or sampling-epistemic strategies at both training and testing time. Ideally, an uncertainty aware-model should achieve high testing reward and high OOD AUC-PR detection score.

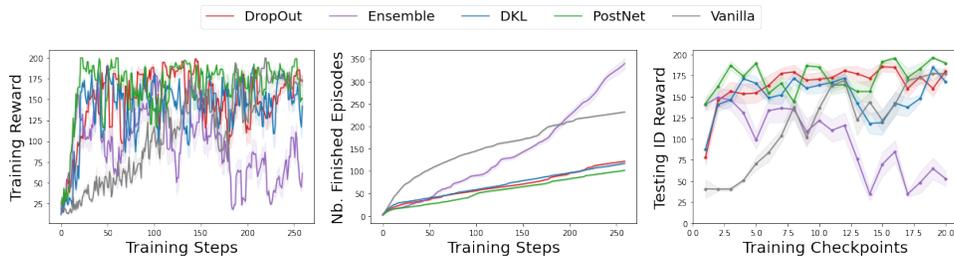


Figure 28: Comparison of the vanilla DQN with the four uncertainty methods performance on Acrobot. All methods use the epsilon-greedy strategy. The vanilla DQN cannot be evaluated on uncertainty tasks.

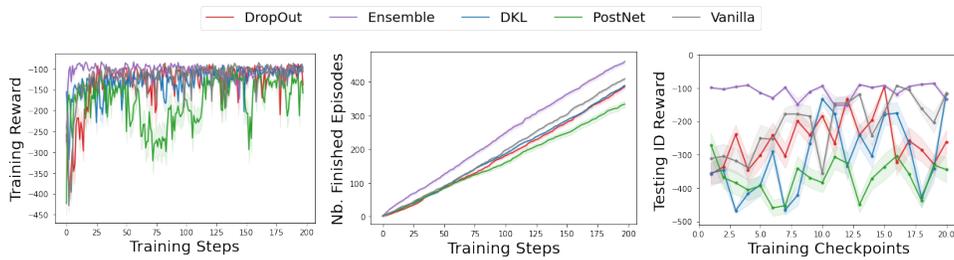


Figure 29: Comparison of the vanilla DQN with the four uncertainty methods performance on Acrobot. All methods use the epsilon-greedy strategy. The vanilla DQN cannot be evaluated on uncertainty tasks.

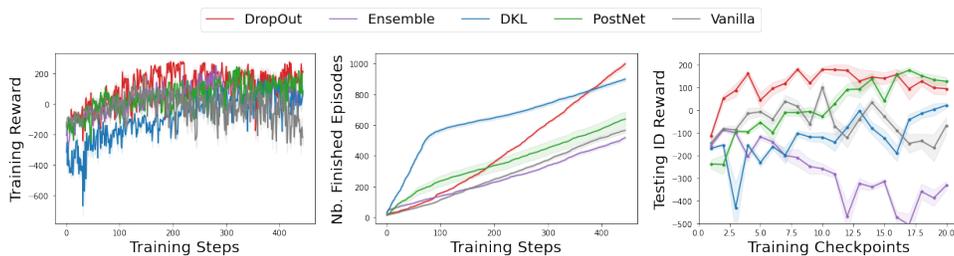


Figure 30: Comparison of the vanilla DQN with the four uncertainty methods performance on LunarLander. All methods use the epsilon-greedy strategy. The vanilla DQN cannot be evaluated on uncertainty tasks.

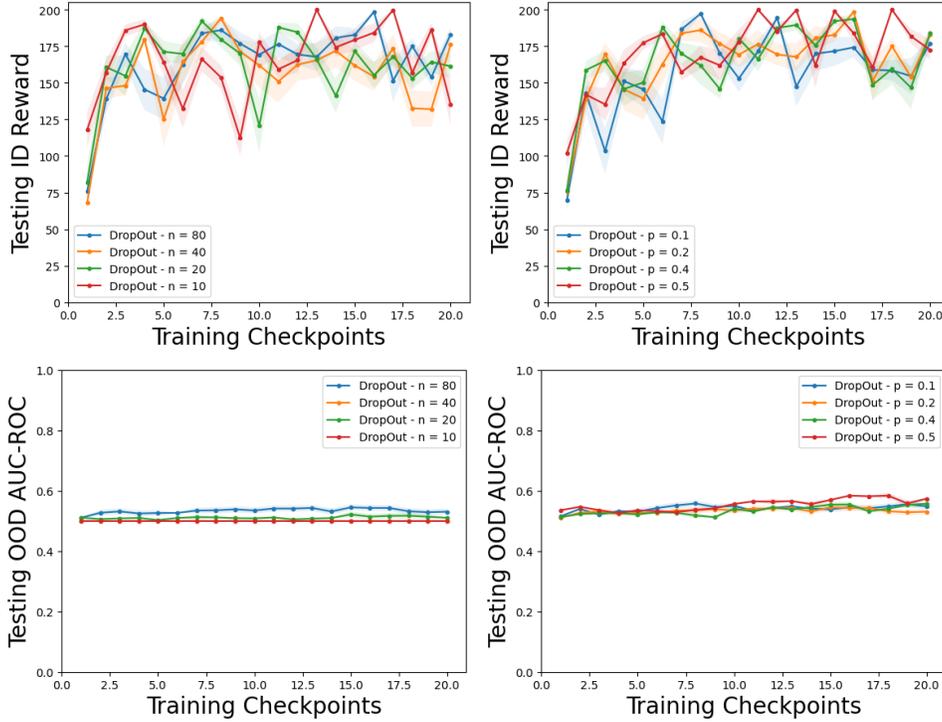


Figure 31: Hyper-parameter study for DropOut w.r.t. the number of samples n and the dropout probability p . Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

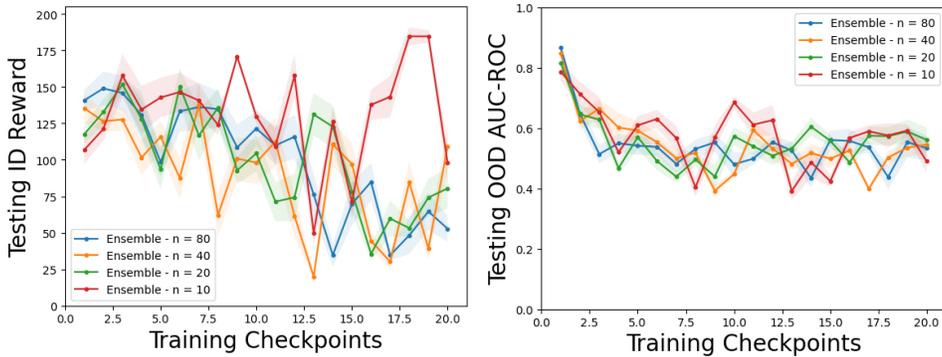


Figure 32: Hyper-parameter study for Ensemble w.r.t. the number of networks n . Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

It facilitates the match between the latent positions output by the encoder and the inducing points. The other hyperparameters consistently show good performances. We show the hyperparameter study for PostNet in fig. 34. Again, the batch norm layer appears to improve the result stability as observed in [17]. It facilitates the match between the latent positions output by the encoder and non-zero density regions learned by the normalizing flows. The other hyperparameters consistently show good performances.

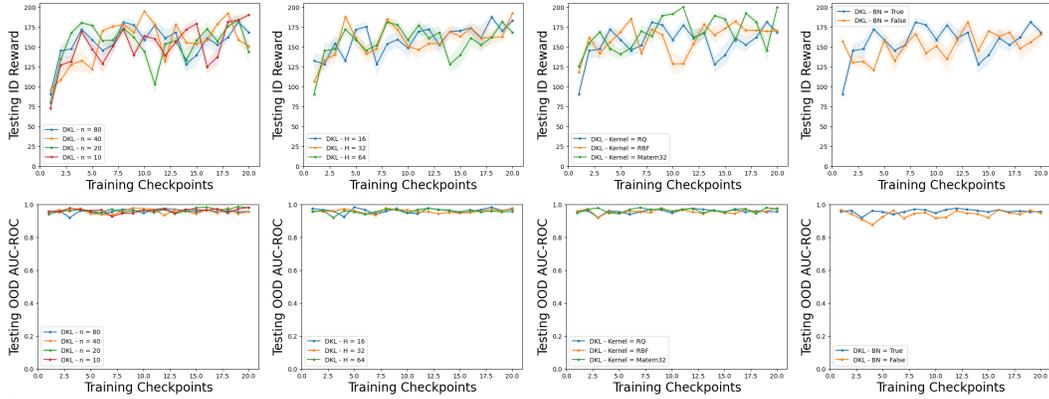


Figure 33: Hyper-parameter study for DKL w.r.t. the number of inducing points n , the latent dimension H , the kernel type and the batch norm layer. Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.

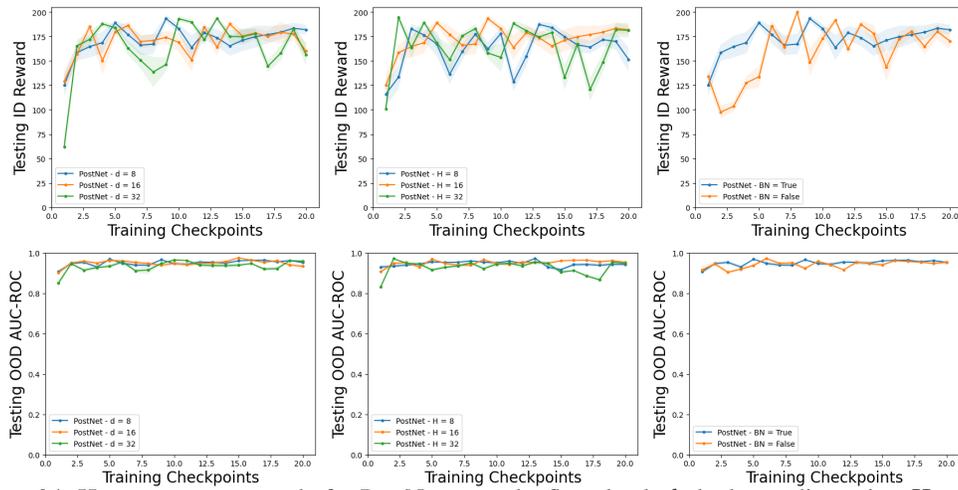


Figure 34: Hyper-parameter study for PostNet w.r.t. the flow depth d , the latent dimension H and the the batch norm layer. Ideally, an uncertainty aware-model should achieve high reward and high OOD detection scores.