

---

# AREPO: Uncertainty-Aware Robot Ensemble Learning Under Extreme Partial Observability

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Real-world applications of vision-based robot learning face two major chal-  
2 lenges: extreme partial observability and effective simulation-to-reality (sim-  
3 to-real) transfer. This paper introduces a robust robot learning framework  
4 that enhances uncertainty awareness to address these challenges. We reinterpret  
5 variational-autoencoder-based visual reinforcement learning (RL) from an  
6 uncertainty-quantification perspective, enabling resilience to high sensory noise and  
7 severe visual occlusions—common in industrial robotic tasks. To further improve  
8 sim-to-real transfer, we propose an uncertainty-aware ensemble RL algorithm.  
9 We validate our methods on a laboratory task designed as a proxy for real-world  
10 industrial applications characterized by harsh environments with low visibility and  
11 physical occlusions. Both simulation and real-world results demonstrate significant  
12 improvements in task accuracy and efficiency over various baselines, highlighting  
13 the benefits of uncertainty-aware robot learning for complex operational contexts.

## 14 1 Introduction

15 In vision-based robot learning, using partial and noisy observations for policy learning presents a  
16 substantial challenge, as a single observational frame often violates the Markov assumption—i.e., that  
17 the current observable state contains all necessary information for future decision-making. To address  
18 this, three main solutions have emerged: The first involves using recurrent neural networks (RNNs)  
19 to encode entire past trajectories [10, 9]. While RNNs theoretically capture comprehensive historical  
20 data, they suffer from slow training times and high computational costs due to the sparse nature of  
21 reinforcement learning (RL) losses, making them less viable for industrial applications. The second  
22 approach constructs a belief state as a statistical summary of past trajectories [14, 8]. However, this  
23 method also incurs high computational costs due to the complexity of continuously updating and  
24 tracking these dynamic states, compounded by the demand for intensive sequence modeling. The  
25 third approach extracts denoised, compressed latent representations from noisy observations using an  
26 autoencoder, effectively revalidating the Markov assumption [11, 35]. This allows for the application  
27 of efficient, memoryless RL methods, reducing computational overhead and better suiting real-world  
28 robotic applications. These autoencoder-based approaches provide a promising compromise between  
29 computational efficiency and the ability to handle complex, noisy data streams.

30 While end-to-end RL allows agents to learn directly from physical interactions with the environment,  
31 industrial robotic tasks often preclude large-scale data collection due to high costs and safety risks.  
32 To mitigate this, RL is frequently integrated into a simulation-to-reality (sim-to-real) framework,  
33 where policies are first developed in simulation before being transferred to real-world settings.

34 Sim-to-real policy transfer is well-established in deep reinforcement learning (DRL), but existing  
35 methods often overlook the challenges of observational uncertainty. Many approaches assume near-  
36 complete observability, rely on precise environmental modeling through domain randomization or

physics-based methods [26, 4], or require continuous adaptation with cheap access to the target domain [21, 30]. While effective in controlled environments, these methods may struggle in the presence of extreme noise and partial observability, where occlusions obscure true state changes. To address this, platform-agnostic transfer methods [19] mitigate task-specific noise by decoupling perception from control, while Kalman filtering integrated with DRL [17, 16] improves robustness in dynamic, noisy tasks. However, both approaches still depend on accurate environmental modeling, which is often impractical in industrial settings due to the complexity of real-world dynamics and constraints on large-scale data acquisition.

To address these challenges, we revisit visual DRL methods based on variational autoencoders (VAE), from an uncertainty-quantification perspective. We explore their connection to performance gaps in sim-to-real transfer, and we introduce a novel uncertainty-aware ensemble DRL framework. Our approach enhances decision-making under extreme noise and partial observability while fostering an uncertainty-based collaborative ensemble mechanism. This mechanism aids in transitioning from potentially inaccurate models during training to effective real-world applications. We achieve this through a unique ensemble learning framework that minimizes deviations among individual policies, encouraging emergent behavior aligned with collective wisdom: an uncertainty-weighted sum of all policies within the ensemble, prioritizing the policy with the least sim-to-real gap.

Our contributions are threefold: (1) We propose uncertainty-aware DRL algorithms based on a reinterpretation of VAE-based visual DRL, leading to improved sampling efficiency (Section 4.1) and enhanced sim-to-real transfer performance (Section 4.2). (2) To validate our approach, we design a laboratory task representative of a wide range of real-world applications, including shotcreting [20], sandblasting [33], and paint spraying [34] (Section 5). (3) We provide a comprehensive survey and reinterpretation of state-of-the-art VAE-based visual DRL from an uncertainty-quantification perspective, addressing the challenges of extreme partial observability and harsh industrial environments (Section 4.1). Our results demonstrate that policies learned in simulation not only generalize effectively to real-world conditions in a zero-shot manner but also outperform traditional model-based planners [2] and DRL baselines [32, 35, 7, 36].

## 2 Related Work

Uncertainty quantification in DRL is critical for providing agents with deeper insights into the learning process. Uncertainty can be categorized into two types: aleatoric and epistemic. Aleatoric uncertainty is inherent to stochastic data and irreducible, whereas epistemic uncertainty arises from the agent’s incomplete understanding of the environment. Most uncertainty quantification methods focus on techniques such as bootstrapping and Monte Carlo (MC) dropout.

Bootstrapped Deep Q-Networks (DQN) [25] introduced epistemic uncertainty quantification through a shared network with multiple heads, where variance in head predictions reflects uncertainty. Extensions penalize highly uncertain states via bootstrapped prior Q-networks [37, 1], reweight Bellman backups for exploration in ensemble frameworks like SUNRISE [18], or estimate epistemic uncertainty using cross-entropy between synthetic and true samples, as in SUMO [27]. While SUMO provides robust estimates, its search-based design incurs high computational costs, limiting real-time applicability.

MC Dropout [6] approximates Bayesian inference through multiple stochastic forward passes. Though effective in supervised learning, its uncertainty quality underperforms variational inference methods [24]. In continuous control tasks, combining MC Dropout with bootstrapped Q-values improves uncertainty estimates [13], but repeated inference introduces scalability challenges for real-time robotics. Hybrid approaches [3] integrating MC Dropout and bootstrapped ensembles disentangle aleatoric and epistemic uncertainties, leading to improved decision-making.

Despite these advancements, the role of uncertainty quantification in sim-to-real transfer remains underexplored. Existing methods emphasize theoretical benchmarks while overlooking partial observability and domain shifts common in real-world robotics. This work bridges that gap by leveraging variational uncertainty quantification and ensemble learning to enhance sim-to-real transfer, addressing epistemic uncertainty in dynamic, noisy environments without the inefficiencies of MC Dropout or search-based methods.

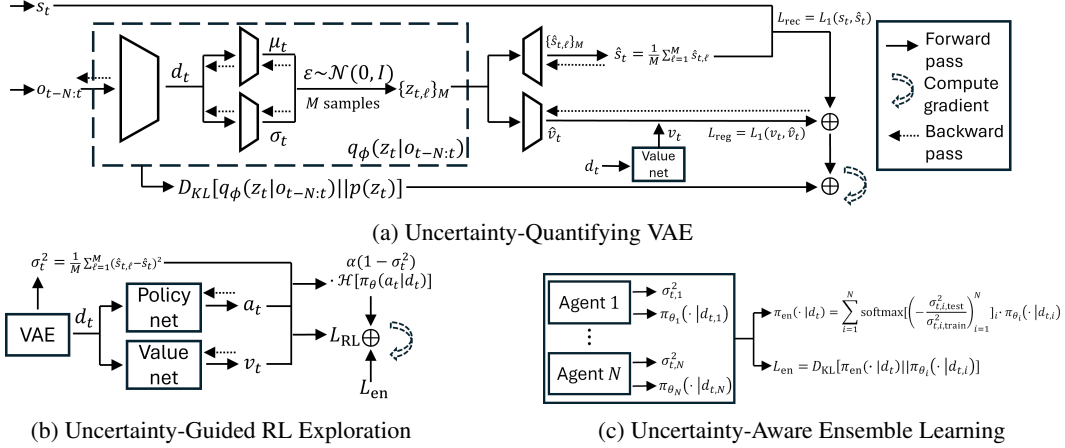


Figure 1: Our algorithm (AREPO) revisits VAE-based RL from an uncertainty quantification perspective. It improves training stability and sampling efficiency with uncertainty-guided exploration and ensemble learning mechanism.

### 3 Preliminaries

In RL, the problem setting is often formulated as an MDP described by the tuple  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ . Here,  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  is the action space,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  is the state transition probability function,  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function, and  $\gamma$  is the discount factor.

The policy  $\pi_\theta$  is optimized to maximize the cumulative reward  $J(\theta)$  of the trajectory  $\tau = (s_t, a_t, r_t, \dots, s_T, a_T, r_T)$ . The cumulative reward is written as

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)}[R(\tau)], \quad (1)$$

where  $R(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}, a_{t'})$ . The policy is generally optimized via a gradient-based method, with  $\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)}[\nabla_\theta \log \pi_\theta(a_t|s_t) \cdot A^{\pi_\theta}(s_t, a_t)]$ . The advantage function  $A^{\pi_\theta}(s_t, a_t)$  is defined as  $A^{\pi_\theta}(s_t, a_t) = Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)}[R(\tau)|s_t, a_t] - \mathbb{E}_{\tau \sim \pi_\theta(\tau)}[R(\tau)|s_t]$ .

Generalized advantage estimation (GAE) is widely used to balance bias and variance, stabilizing training [31]. GAE approximates  $A^{\pi_\theta}(s_t, a_t)$  as  $\hat{A}_{\text{GAE}}^{\pi_\theta}(s_t, a_t) = \sum_{t'=t}^T (\gamma\lambda)^{t'-t} [r(s_{t'}, a_{t'}) + \gamma V_\phi(s_{t'+1}) - V_\phi(s_{t'})]$ , where  $\lambda \in [0, 1]$  is the GAE coefficient, and  $V_\phi$  is the value function estimator learned by minimizing:

$$\mathcal{L}_V(\phi) = \mathbb{E}_{\tau \sim \pi_\theta} [V_\phi(s_{t'}) - R(\tau)]. \quad (2)$$

In this paper, we use proximal policy optimization (PPO) [32], a model-free, on-policy algorithm that constrains policy updates via a clipped surrogate loss:

$$\mathcal{L}_{\text{PPO}}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [\min(w_\theta(s_t, a_t) \cdot \hat{A}_{\text{GAE}}^{\pi_\theta}(s_t, a_t), \text{clip}(w_\theta(s_t, a_t), 1 - \epsilon, 1 + \epsilon) \cdot \hat{A}_{\text{GAE}}^{\pi_\theta}(s_t, a_t))], \quad (3)$$

where  $w_\theta(s_t, a_t) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ , and  $\epsilon$  is the clipping factor.

In real-world robotic applications, agents rely on sensors to obtain observations  $o_t$  of the state  $s_t$ . This complicates RL, as observations are often high-dimensional, noisy, and partially occluded. These challenges have inspired methods that utilize VAEs [11] to learn a denoised, compact latent representation for state inference [12, 5, 23]. A VAE consists of a convolutional encoder  $p_\phi(z_t|o_t)$  that processes noisy, partially occluded observations  $o_t$  and generates a latent representation  $z_t$ , which serves as input to the policy and value functions instead of  $s_t$ . The decoder  $p_\varphi(\hat{s}_t|z_t)$ , structurally mirroring the encoder, reconstructs an estimate of the corresponding state  $\hat{s}_t$  from  $z_t$  using transposed convolutional layers that incrementally upscale  $z_t$ . The VAE is trained by optimizing the following variational lower bound:

$$\mathcal{L}_{\text{VAE}} = -\mathbb{E}_{z \sim q_\phi(z_t|o_t)} [\log p_\varphi(\hat{s}_t|z_t)] + \beta \cdot D_{\text{KL}}[q_\phi(z_t|o_t) || p(z_t)] = \mathcal{L}_{\text{rec}} + \beta \cdot \mathcal{L}_{\text{KL}}, \quad (4)$$

114 where  $q_\phi(z_t|o_t)$  is the variational distribution, and  $p(z_t)$  is the Gaussian prior.

115 Inspired by these works and empirical evidence from [35], we propose a novel VAE-based DRL  
 116 network architecture, where a denoised latent representation  $d_t$  is learned and passed to the policy  
 117 and value networks instead of  $z_t$ . In the next section, we motivate our design choice.

## 118 4 Method

119 In this section, we introduce AREPO (Uncertainty **A**ware **R**obot **E**nsemble Learning under Extreme  
 120 **P**artial **O**bservability). Our approach integrates a denoising VAE for uncertainty-aware state recon-  
 121 struction, an uncertainty-guided RL strategy that prioritizes exploration based on uncertainty, and an  
 122 ensemble method that leverages uncertainty awareness for more effective sim-to-real transfer.

### 123 4.1 Reframing VAE-based Visual RL with Uncertainty Quantification

124 Given a sequence of past observations  $o_{t-N:t}$ , our goal is to learn  $p(d_t|o_{t-N:t})$ , from which we can  
 125 extract a statistically equivalent denoised embedding  $d_t$ . This embedding is theoretically sufficient  
 126 for optimal decision-making because  $p(d_t|o_{t-N:t})$  implicitly contains signals of transition dynamics,  
 127 agent policy, and value functions:

$$\begin{aligned} p(d_t|o_{t-N:t}) &\propto p(o_t|d_t) \cdot p(d_t|o_{t-N:t-1}) \cdot p(o_{t-N:t-1}) = p(d_t)p(o_{t-N:t-1}|d_t) \cdot p(o_t|o_{t-N:t-1}, d_t) \\ &\propto \underbrace{p(d_{t-N}|o_{t-N})}_{\text{encoder}} \prod_{t'=t-N+1}^{t=t} \underbrace{p(o_{t'}|d_{t'})}_{\text{decoder}} \cdot \underbrace{p(d_{t'}|d_{t'-1}, a_{t'-1})}_{\text{dynamics}} \cdot \underbrace{\pi(a_{t'-1}|d_{t'-1})}_{\text{policy}} \cdot \underbrace{q(v_{t'-1}|d_{t'-1})}_{\text{value}}. \end{aligned} \quad (5)$$

128 The observation above is corroborated by positive results from VAE-based DRL [12, 5, 23]: grounding  
 129 RL in an autoencoder’s latent representation, which implicitly encodes the policy, value, and dynamics,  
 130 has been shown to improve sampling efficiency and training stability.

131 However, despite these benefits, previous works based on this method often suffer from suboptimal  
 132 performance in certain scenarios. Multiple empirical workarounds have been proposed in the literature  
 133 to mitigate this issue. These include alternating the training of the VAE and RL components [5, 23],  
 134 using smaller  $\beta$  values in  $\mathcal{L}_{\text{VAE}}$ , or preventing the policy’s gradient from updating other networks  
 135 except itself [35]. In this section, we revisit VAE-based DRL from an uncertainty quantification  
 136 perspective, hypothesize key sources of suboptimality in previous work, and propose an efficient and  
 137 robust uncertainty-guided RL method to address them in a principled manner.

138 Our objective is to guide RL exploration based on the epistemic uncertainty  $p(d_t|o_{t-N:t})$  of the  
 139 neural network that computes  $d_t$  from  $o_{t-N:t}$ . We measure this uncertainty using the variance of  
 140  $p(\hat{s}_t|o_{t-N:t})$ , as it integrates the uncertainty of  $p(d_t|o_{t-N:t})$  when reconstructing the state  $\hat{s}_t$ :

$$p(\hat{s}_t|o_{t-N:t}) = \int p(\hat{s}_t|z_t) \cdot p(z_t|o_{t-N:t}) dz_t = \iint p(\hat{s}_t|z_t) p(d_t|o_{t-N:t}) \mathcal{N}(\mu_t(d_t), \sigma_t(d_t)) dd_t dz_t. \quad (6)$$

141 To compute uncertainty tractably, we follow [24] and use Monte Carlo integration to approximate  
 142 the expectation and variance of the output  $\hat{s}_t$  by drawing  $M$  samples of  $\{z_{t,\ell}\}_{\ell=1}^M$  and obtaining  $M$   
 143 corresponding outputs  $\{\hat{s}_{t,\ell}\}_{\ell=1}^M$ . We use the empirical variance  $\sigma_t^2$  as the uncertainty measure:

$$\sigma_t^2 = \int (\hat{s}_t - \mathbb{E}[\hat{s}_t])(\hat{s}_t - \mathbb{E}[\hat{s}_t])^T p(\hat{s}_t|o_{t-N:t}) d\hat{s}_t \approx \frac{1}{M} \sum_{\ell=1}^M (\hat{s}_{t,\ell} - \mathbb{E}[\hat{s}_t])^2, \quad (7)$$

144 where  $\mathbb{E}[\hat{s}_t] = \int \hat{s}_t \cdot p(\hat{s}_t|o_{t-N:t}) d\hat{s}_t \approx \frac{1}{M} \sum_{\ell=1}^M \hat{s}_{t,\ell}$ .

145 As shown in Fig. 1a, our method differs from previous works in several key aspects:

146 1) We ground the value and policy networks in  $d_t$  instead of  $z_t$ . The reason behind this design  
 147 choice is that we model  $d_t$  as the denoised latent state that is statistically equivalent to  $o_{t-N:t}$ ,  
 148 whereas  $z_t$  is intended solely for epistemic uncertainty quantification of  $p(d_t|o_{t-N:t})$ . Therefore,  
 149  $z_t$  contains undesirable stochasticity and noise for the RL components compared to  $d_t$ , leading to

subpar performance. This reasoning aligns with the empirical findings in [35], where a diminishing  $\beta$  improves performance and stabilizes training.

2) We integrate training schemes from [5, 23] and [35] by combining alternating and joint optimization strategies. [5, 23] stabilize learning by decoupling VAE and RL optimization, while [35] jointly leverages critic gradients and VAE losses to regulate  $z_t$ . To merge these approaches, we introduce a value head  $p(\hat{v}_t|z_t)$  parallel to the VAE decoder, mimicking the critic’s role in guiding representation learning. This design retains the stability of alternating training while benefiting from critic-driven regularization.

3) Our method uses  $\sigma_t^2$  to dynamically adjust the maximum entropy objective first proposed in [38] to effectively enhance sampling efficiency. Equation (1) is replaced with an altered maximum entropy RL objective:

$$\pi_\theta^* = \arg \max_{\pi_\theta} \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}, a_{t'}) + \alpha_1 (1 - \sigma_t^2) \mathcal{H}(\pi(\cdot|d_t)) \right], \quad 0 < \sigma_t^2 < 1. \quad (8)$$

The entropy coefficient in Equation (8) serves two key purposes. First, it draws inspiration from uncertainty-driven exploration methods such as MOPO [37] and SUMO [27], which penalize rewards in high-uncertainty regions of learned transition dynamics. Although our approach is model-free and does not explicitly learn dynamics, its uncertainty estimator implicitly captures all RL components, as shown in Equations (5-7). This ensures that the performance guarantees from [37] remain valid in our setting, maintaining the benefits of uncertainty-aware exploration. Second, we extend the original maximum entropy objective [38] by introducing a dynamic adjustment mechanism. Instead of a fixed entropy weight, our approach applies a soft penalty that discourages policy distributions failing to reduce the agent’s epistemic uncertainty about the environment. This ensures a principled balance between exploration and exploitation. In Section 5.2, we empirically demonstrate that this adaptive entropy regulation significantly improves sample efficiency compared to methods without uncertainty guidance.

## 4.2 Approaching Zero-Shot Policy Transfer with Uncertainty Aware Ensemble Learning

Assuming that we can sample  $N$  MDPs during training and denote the  $i^{\text{th}}$  MDP as  $M_i$ , our goal is to train a policy  $\pi$  that achieves the best expected test-time reward. Rather than learning a globally optimal policy as a generic POMDP solver might, recent advances in DRL often reduce the global policy learning problem into a set of local policy learning problems to scalably solve MDPs. One such example is LEEP [7], which addresses RL generalization by enabling independent learning for each agent and facilitating cross-agent knowledge sharing through a linker function. In LEEP, the optimal global policy  $\pi^*$  is obtained by combining local policies  $\pi_i^*$  optimal in their respective MDPs using a linker function  $\pi^* = f(\{\pi_i^*\}_{i \in [N]})$ , where  $\{\pi_i^*\}_{i \in [N]} = \arg \max_{\pi_1, \dots, \pi_N} \frac{1}{N} \sum_{i=1}^N J_{M_i}(\pi_i) - \alpha_2 \cdot \sum_{i=1}^N \mathbb{E}_{p_{M_i}^{\pi_i}} [D_{\text{KL}}[\pi_i || f(\{\pi_i\}_{i \in [N]})]]$ .

However, LEEP proposes a probability-based linker function that combines policies by selecting the most probable action. This can limit the capacity of ensemble learning, as each agent is likely to prioritize its own policy during training, and the sim-to-real gap is not explicitly optimized. To address this issue, we derive an uncertainty-based generalization gap and use it to assign weights to each agent’s policy contribution in the linker function. This ensures that the combined policy directly optimizes the generalization gap, leading to improved sim-to-real transfer performance.

The generalization gap from [15, 29] can be written as  $\mathbb{E}_{\tau, s \sim p_{M_{\text{train}}}^\pi} [R(\tau)|s] - \mathbb{E}_{\tau, s \sim p_{M_{\text{test}}}^\pi} [R(\tau)|s]$ . Recent work has shown that this gap can be measured by the latent representation deviation  $\Delta d_t$  between training and testing environments, i.e.,  $\|d_{t, \text{test}} - d_{t, \text{train}}\|$  [22]. However, this metric is impractical in real-world robotic policy transfer because it requires  $d_t^{\text{train}}$  and  $d_t^{\text{test}}$  to share the same underlying state  $s_t$ , which is often infeasible to obtain due to noisy and partial observations.

To tractably quantify the generalization gap, we relate  $\Delta d_t$  to the deviation in the reconstructed state  $\Delta \hat{s}_t$ , i.e.,  $\|\hat{s}_{t, \text{test}} - \hat{s}_{t, \text{train}}\|$ . Following [22], we assume Lipschitz continuity in a set of functions that map  $\hat{s}_t$  to  $d_t$  with a constant  $L > 0$ , yielding  $\Delta d_t \leq L \Delta \hat{s}_t$ . Similar to [37], we assume an admissible error estimator with upper bounds  $\sigma_{t, \text{train}}^2$  and  $\sigma_{t, \text{test}}^2$  on the reconstruction errors for both training and testing. Applying the triangle inequality, we obtain:  $\Delta d_t \leq L(\sigma_{t, \text{test}}^2 + \sigma_{t, \text{train}}^2)$ . We define the

generalization gap as a vector  $\delta = (\delta_i)_{i=1}^N$ , where  $\delta_i$  denotes the generalization gap of agent  $i$ . By normalizing  $\Delta d_{t,i}$  by the continuity constant and the estimator error bound, we define:

$$\delta_i = 1 + \frac{\sigma_{t,i,\text{test}}^2}{\sigma_{t,i,\text{train}}^2}. \quad (9)$$

To leverage the derived uncertainty-based generalization gap for ensemble learning, we propose a linker function based on the generalization gap to combine policies as a weighted sum of individual policies, with weights inversely proportional to the generalization gap. This ensures that agents with smaller sim-to-real gaps have a greater influence on the combined policy, while still incorporating knowledge from other agents:

$$\pi_{\text{en}} = \sum_{i=1}^N \text{softmax}(-\delta)_i \cdot \pi_i(\cdot | d_{t,i}). \quad (10)$$

We further compute the KL-divergence between each agent’s policy and the combined policy, serving as an additional loss term during training. This loss encourages coherence across the ensemble, ensuring that individual policies do not diverge significantly from the ensemble consensus:

$$\mathcal{L}_{\text{en}} = D_{\text{KL}}[\pi_{\text{en}} || \pi_i(\cdot | d_{t,i})]. \quad (11)$$

We propose AREPO, a novel ensemble learning architecture that manages individual agent-environment interactions via shared uncertainty quantification in Fig. 1c. Each agent in the ensemble is instantiated with the uncertainty-quantifying VAE and uncertainty-guided RL modules, described (see Fig. 1a and Fig. 1b respectively). We summarize AREPO in Alg. 1 and Alg. 2.

---

**Algorithm 1** AREPO (*training*)

---

```

1: Initialize  $N$  training environments with customized parameterizations, initialize policy  $\theta_i$ , value
    $\eta_i$ , and VAE network  $\phi_i$  for each indexed environment  $E_i, i \in N$ .
2: for each iteration do
3:   for  $i \in N$  do
4:     Collect data  $\{s, o, a, s', o'\}_D$  using  $\pi_i$  in  $E_i$ .
5:     for each mini-batch  $B$  in  $D$  do
6:       Similar to [5] [23], we train the VAE and RL parts in an alternating fashion to stabilize
       training. Activate VAE network  $\phi_i$ .
7:       for  $j \in N, j \neq i$  do
8:         Freeze all networks of agent  $j$ , i.e.,  $\pi_j, \eta_j, \phi_j$ , so agent  $j$  is not updated.
9:       end for
10:      Compute and update VAE network  $\phi_i$  with  $\mathcal{L}_{\text{VAE}}$  computed with Equation (4) and the
       regularization loss  $\mathcal{L}_{\text{reg}} = L_1(v_t, \hat{v}_t)$  in Fig. 1a.
11:      Freeze VAE network  $\phi_i$ . Activate policy and value networks  $\theta_i$  and  $\eta_i$ .
12:      Use Equation (3) to compute  $\mathcal{L}_{\text{ppo}}$  with altered entropy term in Equation (8) and compute
        $\mathcal{L}_V$  with Equation (2).  $\mathcal{L}_{\text{RL}}$  in Fig. 1b is computed using  $\mathcal{L}_{\text{ppo}} + \mathcal{L}_V$ .
13:      Compute  $\sigma_{t,i,\text{train}}^2 = \mathbb{E}_{\sigma_{t,i}^2 \sim S_{i,\text{train}}}[\sigma_{t,i}^2]$ , where  $S_{i,\text{train}}$  is the set of all computed  $\sigma_{t,i}^2$  during
       training. Compute uncertainty measure  $\sigma_{t,i}^2$  with Equation (7) as  $\sigma_{t,i,\text{test}}^2$ . Compute  $\delta_i$ 
       with Equation (9). Compute combined ensemble policy  $\pi_{\text{en}}$  with Equation (10). Compute
       ensemble loss  $\mathcal{L}_{\text{en}}$  with Equation (11).
14:      Update  $\theta_i, \eta_i$  using  $\mathcal{L} = \mathcal{L}_{\text{RL}} + \mathcal{L}_{\text{en}}$ . Freeze policy and value networks  $\theta_i$  and  $\eta_i$ .
15:     end for
16:   end for
17: end for

```

---

## 213 5 Experiments

214 This section presents the experimental evaluation of AREPO, addressing two primary questions:

215 1) Can AREPO achieve improved sampling efficiency and stability compared to uncertainty-unaware  
 216 DRL baselines under extreme partial observability?

---

**Algorithm 2** AREPO (*inference*)

---

- 1: Load  $N$  trained agents, initialize policy  $\theta_i$ , value  $\eta_i$ , VAE network  $\phi_i$ , and epistemic training uncertainty  $\sigma_{t,i,\text{train}}^2$  for each agent  $i, i \in N$ .
  - 2: **while** the task is not done **do**
  - 3:   Receive observation  $o_{t-N:t}$ .
  - 4:   **for** each agent  $i \in N$  **do**
  - 5:     Compute denoised latent state  $d_{t_i}$  using VAE network  $\phi_i$  and  $o_{t-N:t}$ .
  - 6:     Compute the corresponding epistemic uncertainty  $\sigma_{t,i}^2$  with Equation (7) as  $\sigma_{t,i,\text{test}}^2$ .
  - 7:     Compute  $\delta_i$  with Equation (9). Compute policy  $\pi_i(\cdot|d_{t,i})$  using policy network  $\theta_i$  and  $d_{t,i}$ .
  - 8:   **end for**
  - 9:   Compute combined policy  $\pi_{\text{en}}$  with Equation (10). Return action  $a \sim \pi_{\text{en}}$ .
  - 10: **end while**
- 

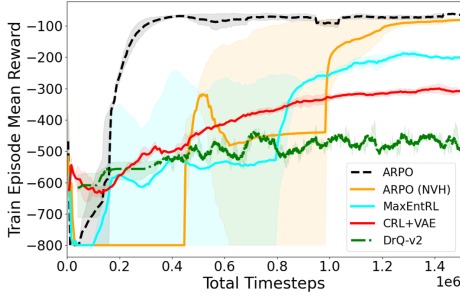


Figure 2: ARPO achieves the best sampling efficiency and training stability under extreme partial observability compared to DRL methods without uncertainty awareness.

Table 1: In the experiment on sampling efficiency and training stability, ARPO obtains the best surface quality and most efficient material usage compared to uncertainty-unaware baselines.

Metrics	$R_{\text{rms}}$ (mm)	$R_t$ (mm)	$r_{\text{av}}$ (%)	$t_{\text{avg}}$ (ms)
ARPO	<b><math>0.6 \pm 0.2</math></b>	<b><math>3.9 \pm 1.4</math></b>	<b><math>23.8 \pm 0.01</math></b>	$17.2 \pm 1.4$
CRL+VAE	$1.7 \pm 0.9$	$8.3 \pm 4.9$	$32.5 \pm 0.02$	$16.6 \pm 2.3$
MaxEntRL	$1.3 \pm 0.6$	$5.5 \pm 2.5$	$33.8 \pm 0.01$	<b><math>7.7 \pm 0.5</math></b>
MPC	$2.8 \pm 0.4$	$13.5 \pm 1.6$	$31.3 \pm 0.05$	$59.2 \pm 2.5$
Vanilla	$10.2 \pm 0.4$	$40.6 \pm 1.5$	$53.8 \pm 0.02$	$50.1 \pm 1.2$

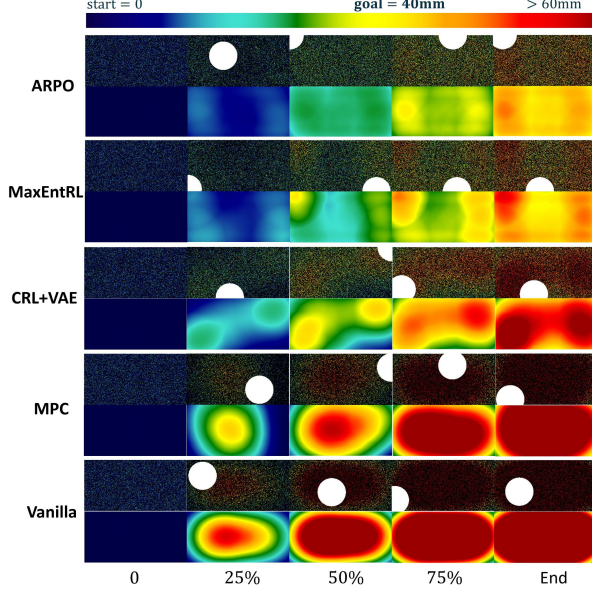


Figure 3: For each method, the first row displays the observed heightmaps used as input to the models, where the white areas are unobservable to agents. The second row shows corresponding states. The color shift in the heightmaps indicates task progress.

217 2) Can AREPO generalize to novel scenarios characterized by temporally-correlated dynamics and  
218 varying levels of spatially-correlated partial observability?

219 To investigate these questions, we evaluate AREPO on a task that involves guiding a tool to sprinkle  
220 material onto a target surface to a specified thickness (See Fig. 4, displayed on page 8 to conve-  
221 niently appear next to its sibling Fig. 5). This task serves as a proxy for industrial applications  
222 like shotcreting [20], sandblasting [33], and paint spraying [34], where robots must operate under  
223 severe visual occlusions and sensory noise. By contrast, standard RL benchmarks used in previous  
224 works (e.g., Control Suite [9, 35], Mujoco [36, 30], or Atari [18, 14]) fail to capture temporally  
225 and spatially correlated occlusions and noise common in real-world industrial settings. To further

Table 2: In the experiment on zero-shot sim-to-sim transfer, AREPO remains robust under varying levels of extreme partial observability, achieving better surface quality and robustness than LEEP due to the uncertainty weighted ensemble policy. DR generally achieves the worst performance, indicating a failed sim-to-sim transfer to tasks under extreme partial observability.

Scenarios	Scenario 1: 65% Random Occlusion + 15% Plume				Scenario 2: 80% Random Occlusion + 50% Plume				Scenario 3: 95% Random Occlusion + 85% Plume			
Metrics	$R_{\text{rms}}$ (mm)	$R_t$ (mm)	$r_{\text{av}}$ (%)	$t_{\text{avg}}$ (ms)	$R_{\text{rms}}$ (mm)	$R_t$ (mm)	$r_{\text{av}}$ (%)	$t_{\text{avg}}$ (ms)	$R_{\text{rms}}$ (mm)	$R_t$ (mm)	$r_{\text{av}}$ (%)	$t_{\text{avg}}$ (ms)
AREPO	<b><math>1.6 \pm 0.3</math></b>	<b><math>8.1 \pm 1.3</math></b>	$30.0 \pm 0.01$	$32.3 \pm 2.2$	<b><math>1.7 \pm 0.5</math></b>	<b><math>8.6 \pm 1.9</math></b>	$30.0 \pm 0.01$	$32.1 \pm 3.2$	<b><math>1.7 \pm 0.3</math></b>	<b><math>8.8 \pm 1.4</math></b>	<b><math>31.3 \pm 0.01</math></b>	$30.7 \pm 2.7$
LEEP	$3.6 \pm 0.6$	$16.6 \pm 2.8$	$42.5 \pm 0.03$	$28.1 \pm 3.6$	$3.9 \pm 0.6$	$17.9 \pm 2.5$	$40.2 \pm 0.03$	$28.9 \pm 3.1$	$4.6 \pm 0.5$	$20.7 \pm 2.0$	$47.5 \pm 0.02$	$27.2 \pm 2.5$
DR	$4.3 \pm 2.6$	$18.4 \pm 4.0$	$58.8 \pm 0.07$	<b><math>17.5 \pm 1.2</math></b>	$4.5 \pm 3.4$	$18.5 \pm 14.0$	$62.5 \pm 0.06$	<b><math>17.4 \pm 1.3</math></b>	$6.0 \pm 3.4$	$25.7 \pm 10.1$	$71.3 \pm 0.01$	<b><math>17.3 \pm 1.5</math></b>
MPC	$2.8 \pm 0.4$	$13.5 \pm 1.6$	$31.3 \pm 0.05$	$59.2 \pm 2.5$	$3.1 \pm 0.5$	$13.6 \pm 1.9$	$32.5 \pm 0.04$	$59.3 \pm 3.2$	$3.9 \pm 0.5$	$17.1 \pm 1.7$	$37.5 \pm 0.03$	$59.8 \pm 2.4$
Vanilla	$8.9 \pm 0.6$	$35.7 \pm 2.36$	$48.8 \pm 0.04$	$50.2 \pm 1.2$	$10.1 \pm 0.7$	$40.3 \pm 2.5$	$52.5 \pm 0.03$	$50.1 \pm 1.2$	$10.2 \pm 0.4$	$40.6 \pm 1.5$	$53.8 \pm 0.02$	$50.1 \pm 1.2$

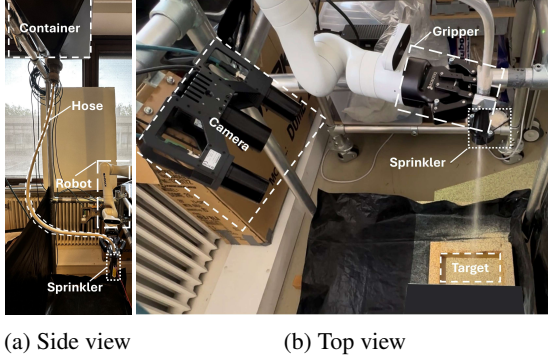


Figure 4: Our *sand sprinkling* testbed serves as a proxy for a variety of industrial applications characterized by significant visual occlusions and noise. Fig. 4a: Sand stored in an overhead container is transmitted through a hose to a sprinkler positioned over a target (a rectangular bin). Sprinkling generates plumes that create visual occlusions. Fig. 4b: The sprinkler is held by a robotic arm which is guided using heightmaps of the target, derived from images captured by an overhead stereo camera.

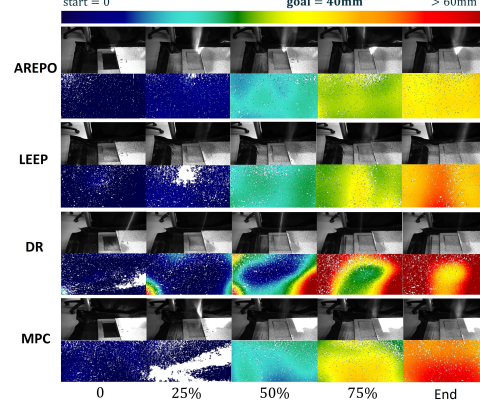


Figure 5: For each method, the first row displays the grayscale images of the scenario. The second row shows their corresponding heightmaps of the target surface. In the sim-to-real experiment, AREPO is most robust to partial observability caused by dust and sensory noise (unobservable white parts of the heightmaps) and achieves best surface quality compared to other methods.

Table 3: In zero-shot sim-to-real transfer to the laboratory testbed, AREPO achieves the best surface quality compared to baselines without uncertainty awareness.

Metrics	$R_{rms}$ (mm)	$R_t$ (mm)	$r_{wv}$ (%)	$t_{avg}$ (ms)
AREPO	$1.1 \pm 0.5$	$7.6 \pm 1.5$	$24.2 \pm 2.1$	$30.2 \pm 3.2$
LEEP	$3.3 \pm 0.6$	$17.1 \pm 2.8$	$30.2 \pm 2.5$	$29.2 \pm 2.1$
DR	$4.3 \pm 1.8$	$21.2 \pm 4.3$	$45.7 \pm 4.3$	$16.7 \pm 1.5$
MPC	$3.7 \pm 0.3$	$18.7 \pm 2.2$	$32.4 \pm 3.1$	$61.1 \pm 3.7$

illustrate the task’s difficulty and relevance, we evaluate a Model Predictive Control (MPC) [2] planner, and a vanilla planner that operates by always guiding the tool to the position with the least material deposition alongside DRL methods. As is shown in Sections 5.3 and 5.4, the performance deteriorates significantly for all methods except AREPO in scenarios with increasing noise and occlusion, highlighting the difficulty of the task and its suitability for evaluating policies in noisy, partially observable environments.

To address **Question 1**, we evaluate ARPO (the non-ensemble version of AREPO) in simulation against several baselines: DrQ-v2 (encoder-only DRL) [36], MaxEntRL (PPO + maximum entropy objective [38]+VAE), CRL+VAE (MaxEntRL + critic-guided representation learning [35]), ARPO-NVH (CRL+VAE + uncertainty-guided exploration), and ARPO (ARPO-NVH + value head). In the synthetic experiment addressing **Question 1**, we simulate spatially correlated noise, with severe occlusions concentrated near the contact point where the material meets the surface. Portions of the target surface are also obscured, simulating real-world conditions where dust, debris, and airborne particles reduce visibility and complicate task execution. The material flow rate is modeled as temporally correlated and episodic, with adjustable episode durations to reflect fluctuating pressures or discharge rates typical in real-world operations. Fig. 3 shows a visual representation of the data generated by our simulator. To address **Question 2**, We evaluate MPC, DR (vanilla domain randomization) [26], and LEEP (ensemble-based transfer learning method) [7] against AREPO (LEEP + uncertainty-weighted linker function in Equation 9). All DRL methods use ARPO as their agents for a fair comparison. Experiments addressing **Question 2** consist of two training scenarios, both carried out in simulation: one with a constant material flow rate and another with a extremely varied flow rate ranging from 10% to 1000% of the first rate. Both scenarios feature the same spatially correlated partial observability as in the previous experiment. After training, the evaluation of **Question 2** is conducted in two phases. In the first phase, zero-shot tests are performed in simulation under varying degrees of extreme occlusion. This phase evaluates how AREPO and the baseline models adapt to severe partial observability in a controlled simulation. In the second phase, the learned policies are

transferred zero-shot to the laboratory testbed. For a detailed layout of the laboratory setup, see Fig. 4a and Fig. 4b.

## 5.1 Technical Details

All models in the experiment receive input in the form of 2D heightmaps, representing the material deposition status on the target surface. The output of each model is a 2D velocity vector that controls the nozzle’s position in a plane parallel to the target, determining where material is applied. To account for variability during evaluation, the performance metrics are averaged over 20 independent validation trials. All DRL models are implemented using Stable-Baselines3 [28], and trained for 1.6 million timesteps before evaluation. Full details of model architectures and hyperparameters can be found in our publicly available source code at <https://gitlab.kuleuven.be/detry-lab/public/arepo>. The performance of the agents are compared with 4 metrics: i) *root-mean-square roughness*  $R_{\text{rms}}$ : the square root of the mean of the squares of the deviations of the surface height values from the mean surface height, ii) *peak-to-valley roughness*  $R_t$ : the difference in height between the highest point and the lowest point on a surface, iii) *waste volume ratio*  $r_{\text{wv}}$ : the ratio between the wasted volume and the desired volume to be fulfilled. The wasted volume is defined as the material volume that has been sprayed outside the target surface or that exceeds the target thickness. iv) *average inference time*  $t_{\text{avg}}$ : the average time the agent takes to compute  $a_t$  given  $o_{t-N:t}$ .

## 5.2 Simulated Experiment on Sampling Efficiency and Training Stability

Fig. 2 reveals several key insights. First, VAE-based DRL enables more robust representation learning as compared to DrQ-v2 that relies on an encoder-only DRL. Second, CRL+VAE, which augments learning with critic loss, is more stable than MaxEntRL but is more prone to local optimality. Finally, while ARPO-NVH outperforms CRL+VAE by avoiding local optimality using uncertainty-guided exploration, its stability is lower. The inclusion of a value head in ARPO stabilizes training further, highlighting the importance of the value head in balancing exploration and stability. Additionally, Fig. 3 demonstrates ARPO’s superior surface quality and minimal material waste, as supported by Tab. 1.

## 5.3 Simulated Experiment on Zero-Shot Sim-to-Sim Transfer

Tab. 2 shows that AREPO outperforms LEEP across varying levels of partial observability. This success is attributed to the uncertainty-weighted ensemble policies, which enable AREPO to adapt effectively to occlusion and noise in the environment. In contrast, DR and MPC agents exhibit significant performance deterioration as occlusion increases, indicating failure in sim-to-sim transfer under extreme partial observability. While the vanilla planner shows robust performance due to its simple rules, this simplicity also limits its effectiveness in handling more complex tasks.

## 5.4 Real-World Experiment on Zero-Shot Sim-to-Real Transfer

As illustrated in Fig. 5, AREPO demonstrates superior surface quality, achieving a more homogeneous material application across the target surface compared to other methods. Tab. 3 further quantifies this performance, showing that DR and MPC struggle to generalize effectively, with DR producing the poorest surface quality. Both LEEP and AREPO outperform DR and MPC, highlighting the value of ensemble DRL approaches. AREPO’s superiority over LEEP is attributed to its uncertainty-weighted combined policy, which dynamically adjusts based on real-time uncertainty estimations, resulting in more robust material application.

## 6 Conclusions

This paper introduces AREPO, an ensemble DRL framework that enhances robustness under extreme partial observability by incorporating uncertainty quantification through a VAE-based approach, thereby improving sampling efficiency during exploration and training stability via an additional value prediction head. This uncertainty estimation bridges the sim-to-real generalization gap, enabling a more adaptive ensemble mechanism that leverages real-time uncertainty. Empirical results demonstrate AREPO’s superior learning efficiency and policy robustness compared to traditional methods.

## References

- [1] C. Bai and *et al.* Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In *International Conference on Learning Representations*, 2021. Oct. 2021.
- [2] E. F. Camacho and C. Bordons. *Model Predictive Control*. Springer-Verlag London Limited, 2nd edition, 2007.
- [3] B. Charpentier, R. Senanayake, M. Kochenderfer, and S. Günnemann. Disentangling epistemic and aleatoric uncertainty in reinforcement learning. *arXiv preprint arXiv:2206.01558*, 2022.
- [4] F. Djeumou and *et al.* One model to drift them all: Physics-informed conditional diffusion model for driving at the limits. In *Conference on Robot Learning*, pages 123–134, 2024.
- [5] C. Finn and *et al.* Learning visual feature spaces for robotic manipulation with deep spatial autoencoders. *arXiv preprint arXiv:1509.06113*, 2015.
- [6] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, pages 1050–1059, 2016.
- [7] 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*, volume 34, pages 25502–25515, 2021.
- [8] K. Gregor, G. Papamakarios, F. Besse, L. Buesing, and T. Weber. Temporal difference variational auto-encoder. *arXiv preprint arXiv:1806.03107*, 2018.
- [9] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [10] M. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In *2015 AAAI Fall Symposium Series*, 2015.
- [11] I. Higgins and *et al.*  $\beta$ -vae: Learning basic visual concepts with a constrained variational framework. In *ICLR Poster*, volume 3, 2017.
- [12] I. Higgins and *et al.* Darla: Improving zero-shot transfer in reinforcement learning. In *International Conference on Machine Learning*, pages 1480–1490, 2017.
- [13] T. Hiraoka, T. Imagawa, T. Hashimoto, T. Onishi, and Y. Tsuruoka. Dropout q-functions for doubly efficient reinforcement learning. *arXiv preprint arXiv:2110.02034*, 2021.
- [14] M. Igl, L. Zintgraf, T. A. Le, F. Wood, and S. Whiteson. Deep variational reinforcement learning for pomdps. In *International Conference on Machine Learning*, pages 2117–2126, 2018.
- [15] M. Jiang, E. Grefenstette, and T. Rocktäschel. Prioritized level replay. In *International Conference on Machine Learning*, pages 4940–4950, 2021.
- [16] E. Kaufmann, L. Bauersfeld, A. Loquercio, and *et al.* Champion-level drone racing using deep reinforcement learning. *Nature*, 620, 2023.
- [17] H. Lee, M. Kim, S. Park, and D. Kim. Domain randomization and improved simulation accuracy for ct robots. *IEEE Robotics and Automation Letters*, 8(3):2345–2352, 2023.
- [18] K. Lee, M. Laskin, A. Srinivas, and P. Abbeel. Sunrise: A simple unified framework for ensemble learning in deep reinforcement learning. In *International Conference on Machine Learning*, pages 6131–6141, 2021.
- [19] D. Li and O. Okhrin. A platform-agnostic deep reinforcement learning framework for effective sim2real transfer towards autonomous driving. *Nature Communications*, 3, 2024.
- [20] B. Lu, M. Li, T. N. Wong, and S. Qian. Effect of printing parameters on material distribution in spray-based 3d concrete printing (s-3dcp). *Automation in Construction*, 124, 2021.

- [21] J. Lyu and *et al.* Cross-domain policy adaptation by capturing representation mismatch. *arXiv preprint arXiv:2405.15369*, 2024.
- [22] J. Lyu, L. Wan, X. Li, and Z. Lu. Understanding what affects generalization gap in visual reinforcement learning: Theory and empirical evidence. *arXiv preprint arXiv:2402.02701*, 2024.
- [23] A. V. Nair and *et al.* Visual reinforcement learning with imagined goals. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [24] I. Oleksienko, D. T. Tran, and A. Iosifidis. Variational neural networks. *Procedia Computer Science*, 222:104–113, 2023.
- [25] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [26] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *International Conference on Robotics and Automation*, pages 3803–3810, 2018.
- [27] Z. Qiao, J. Lyu, K. Jiao, Q. Liu, and X. Li. Sumo: Search-based uncertainty estimation for model-based offline reinforcement learning. *arXiv preprint arXiv:2408.12970*, 2024.
- [28] A. Raffin and *et al.* Stable-baselines3: Reliable reinforcement learning implementations. <https://github.com/DLR-RM/stable-baselines3>, 2021.
- [29] R. Raileanu and R. Fergus. Decoupling value and policy for generalization in reinforcement learning. In *International Conference on Machine Learning*, pages 8787–8798, 2021.
- [30] A. Rusu, M. Večerík, T. Rothörl, N. Heess, R. Pascanu, and R. Hadsell. Sim-to-real robot learning from pixels with progressive nets. In *Conference on Robot Learning*, pages 262–270, 2017.
- [31] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [32] J. Schulman and *et al.* Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [33] D. J. Taljaard, J. Fourie, and C. P. Kloppers. The effect of sandblasting and bead blasting on the surface finish of dry electrolyte polishing of laser powder bed fusion parts. *Journal of Materials Engineering and Performance*, 32:2050–2061, 2023.
- [34] G. Tanaka, Y. Takahashi, and H. Iwata. High precision paint deposition modeling considering variable posture of spray painting robot. In *International Conference on Robotics and Automation*, pages 2542–2548, 2024.
- [35] D. Yarats and *et al.* Improving sample efficiency in model-free reinforcement learning from images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 2021.
- [36] D. Yarats and *et al.* Mastering visual continuous control: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021.
- [37] T. Yu and *et al.* Mopo: Model-based offline policy optimization. In *Advances in Neural Information Processing Systems*, volume 33, pages 14129–14142, 2020.
- [38] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. In *Proceedings of the 23rd National Conference on Artificial Intelligence*, 2008.