RISE: Robust Imitation through Stochastic Encodings

Mumuksh Tayal, Manan Tayal and Ravi Prakash

Abstract-Ensuring safety in robotic systems remains a fundamental challenge, especially when deploying offline policylearning methods such as imitation learning in dynamic environments. Traditional behavior cloning (BC) often fails to generalize when deployed without fine-tuning because it does not account for disturbances in observations that arises in realworld, changing environments. To address this limitation, we propose RISE (Robust Imitation through Stochastic Encodings), a novel imitation-learning framework that explicitly addresses erroneous measurements of environment parameters into policy learning via a variational latent representation. Our framework encodes parameters such as obstacle state, orientation, and velocity into a smooth variational latent space to improve test time generalization. This enables an offline-trained policy to produce actions that are more robust to perceptual noise and environment uncertainty. We validate our approach on two robotic platforms, an autonomous ground vehicle and a Franka Emika Panda manipulator and demonstrate improved safety robustness while maintaining goal-reaching performance compared to baseline methods.

I. INTRODUCTION

As autonomous robots become increasingly integrated into real-world applications, ensuring safe and highperformance control remains a fundamental challenge. To address safety constraints in such scenarios, various optimalcontrol strategies have been explored, including Constrained Model Predictive Control (MPC) [1], Hamilton–Jacobi (HJ) reachability-based methods [2], [3], and Control Barrier Functions (CBFs) [4], [5]. While these approaches provide formal safety guarantees, they typically rely on explicit models of system and environment dynamics, which are often difficult to obtain in real-world settings. Moreover, many of these frameworks assume the ability to perform consistent online rollouts, which raises concerns about the feasibility and safety of conducting unsafe rollouts during training. To address this, several works advocate offline learning approaches [6], [7] that leverage pre-recorded datasets to avoid repeated unsafe rollouts.

At times, safety concerns also extend to demonstration data, where recording unsafe demonstrations may be infeasible. In such cases, Imitation Learning (IL) is a promising approach for training control policies from safe expert demonstrations, particularly when system dynamics are partially known or difficult to model. However, traditional IL methods such as behavioral cloning (BC) often struggle to generalize beyond the training distribution, resulting in degraded performance at deployment, especially when measurement

All the authors belong to Cyber Physical Systems, Indian Institute of Science (IISc), Bengaluru. {mumukshtayal, manantayal, ravipr}@iisc.ac.in.

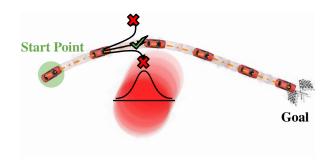


Fig. 1: Behavior Cloning (BC) struggles to generalize for noisy obstacle readings in real-world and often collides with dynamic obstacles, whereas **RISE** picks safer actions without deviating from the original trajectory while it navigates through dynamic obstacles.

equipment provides only rough estimates of obstacle positions and the surrounding environment. Several techniques propose adding noise during demonstration collection [8], [9] to encourage robustness, but this is not always possible, particularly when the demonstration data is pre-recorded.

Attempts have been made to combine methods such as Constrained MPC and CBFs with IL [10], [11] to improve constraint satisfaction in robotic systems. For example, HJ reachability–based imitation learning [9] enforces control constraints by computing forward or backward reachable sets to guarantee safety. However, these methods are computationally expensive and do not scale well to high-dimensional robotic systems due to the curse of dimensionality. Similarly, some CBF-based IL approaches [12] do not explicitly enforce control bounds and often assume unlimited control authority, which is impractical for real systems.

In many practical applications, key safety-related cues (e.g., obstacle position, velocity, and geometry) are available from onboard (or even precise coordinates through motion capture), and task specifications such as goal positions are often pre-defined and may vary across deployments. Although these structured environmental cues are available, their precise and accurate state coordinates and velocity may not be obtainable in real time; instead, only rough measured estimates are typically available.

To address these challenges, we propose RISE, a novel framework that accounts for noisy measurements of obstacle data by conditioning the policy on a learned probabilistic latent space. Specifically, we build on Goal-Conditioned Imitation Learning (GCIL) [13], [14], and augment it by integrating the variational encoder to accommodate noisy, yet safety-critical environmental factors into a structured latent representation. This enables a more realistic interpolation

between environment parameters, thus, improving task adaptability for unseen intermediate datapoints while leveraging the structured perturbation data directly from the latent space. This improves inculcating inherent awareness of safety from the provided data, which is generally hard to achieve with behavior cloning alone.

To summarize, the key contributions of our paper are:

- Unlike CBF and HJ reachability approaches, which
 require exact system dynamics or an accurate environment model, RISE operates in real-world settings where
 dynamics are unavailable or imprecise.
- We train a variational autoencoder (VAE) to predict a
 probabilistic distribution over obstacle states, thereby
 capturing uncertainty from noisy observations. Conditioning the policy on this distribution yields more robust
 behaviors that avoid likely obstacle locations.
- Few safety-critical policy-learning frameworks (including many CBF-based methods) explicitly consider physical actuation limits. Our approach incorporates actuation constraints and learns policies that respect those limits while behaving conservatively near obstacles.
- We validate the framework in simulated robotic environments and on hardware. Comparative analyses against baselines such as PCIL and C-PPO [15], [16] show improved safety with maintained goal-reaching performance.

II. PRELIMINARIES

A. Safe Imitation Learning

Most regular Imitation learning frameworks train policies that map observations to actions using expert demonstrations, in cases where dynamics of the system is unknown or very complex. Various IL approaches, including behavioral cloning (BC) [17], DAgger [18], and inverse reinforcement learning (IRL) [19] exist, however, even though they are derived from the demonstrations of a superior policy, they lack the capability to learn safety-aware policies due to lack of true reward signals. Safe Imitation Learning frameworks, on the other hand, either try to remain within in-distribution region as demonstrated by the expert demonstrations [10], [20], thus, minimizing the risk of violating safety, or they add an adversarial noise to the demonstrations while recording them [8], [9] to inherently learn a more robust policy. But it is important to also consider that these approaches have their shortcomings, either they are overly conservative or they require a specific dataset to train, both of which may not be acceptable at all times.

B. Parameter-Conditioned Imitation Learning

It is a subdomain of Imitation Learning, where each demonstration data-point is augmented with one or more parameters (e.g., goal state [13]), hence, seeking to obtain the indicator reward for the task that the demonstration was provided for. The conditioning parameter contains information that a learning method can leverage to disambiguate demonstrations. Parameters such as goal-states have also extended the domain of reinforcement learning through Goal

Conditioned Reinforcement Learning (GCRL) [21], where the agent is not provided expert demonstrations but reward signals instead. Typically these reward signals are difficult to define, especially for complex tasks and environments, providing demonstrations is often a more natural option in such situations. Additionally, the policy rollouts required by GCRL are often expensive in real-world settings.

C. Variational AutoEncoder (VAE)

Variational AutoEncoders (VAEs) [22] are generative models that learn a probabilistic latent space representation of data. A VAE consists of an encoder and a decoder component, both of which are connected to each other using the reparameterization trick (as referred in [22]).

The objective function of a VAE is to maximize the Evidence Lower Bound (ELBO):

$$\mathcal{ELBO}(x) = \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - D_{KL}(q_{\phi}(z|x)||p(z)), \tag{1}$$

where the first term maximizes the likelihood of reconstructions, and the second term regularizes the latent space by minimizing the Kullback-Leibler (KL) divergence between the approximate posterior and a prior distribution p(z).

VAEs have been widely adopted for learning structured representations, denoising, and improving generalization in downstream tasks, making them a valuable tool for enhancing imitation learning in RL [23], thus, making them an ideal choice for an application like ours.

III. METHODOLOGY

In this section, we present the framework that learns a latent unsafe region distribution for given noisy obstacle perception signals to enable robust imitation learning. The method first encodes raw measured parameters (e.g., obstacle position, obstacle velocity, and obstacle radius as in the cases demonstrated) into a structured latent variable, and then conditions an imitation policy on the current state, goal region and this derived latent distribution.

A. Problem Formulation

Consider a robotic system with state space \mathcal{S} , action space \mathcal{A} , and a set of safety parameters \mathcal{C} . The safety parameters $c \in \mathcal{C}$ represent critical environmental features such as obstacle positions, velocities, and geometries. The objective is to learn a policy $\pi: \mathcal{S} \times \mathcal{C} \to \mathcal{A}$ that maps states and safety parameters to actions while maintaining safety constraints and accomplishing the desired task.

B. Latent Unsafe Region

Let $c \in \mathbb{R}^{d_c}$ denote the raw measured parameters. We employ a variational encoder network $E(\cdot)$ to embed c into a latent normal distribution. To train our VAE-style encoder using ELBO, we model z using the reparameterization trick:

$$z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I),$$
 (2)

where μ is the mean and σ is the deviation, which are the outputs of the encoder network. Thence derived latent

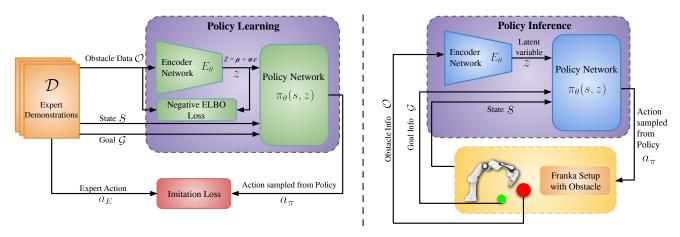


Fig. 2: **Architecture:** (*Left*) During training, expert demonstration trajectories, augmented with obstacle parameters, and goal coordinates, are used to learn the Variational Encoder to produce a latent representation (z), which is combined with agent's state to drive training of the Policy Network. (*Right*) At inference, trained architecture deploys the Policy Network on a real agent, generating actions for safe navigation in dynamic environments.

variable z is passed into the decoder and the VAE is then learnt using the negative ELBO loss.

The encoder architecture consists of fully connected layers with ReLU activations, culminating in parallel output layers for μ and σ .

C. Behavior Policy

The policy $\pi(s,g,z)$ maps the current state $s\in\mathbb{R}^{d_s}$, the goal g and the latent safety variable z to an action $a\in\mathbb{R}^{d_a}$:

$$a = \pi(s, g, z). \tag{3}$$

Unlike traditional behavior cloning approaches that directly map states to actions, our policy also leverages the structured latent representation of measured obstacle parameters. By sampling from the VAE posterior during training, the policy effectively sees multiple plausible obstacle hypotheses (a form of virtual data augmentation), which improves robustness to perceptual variation. The policy network uses a 2-layered fully connected Neural Network with 128 neurons in the hidden layer with ReLU activations. The state vector \boldsymbol{s} , goal vector \boldsymbol{g} and latent vector \boldsymbol{z} are concatenated and passed through these layers to produce the action output.

Algorithm 1 summarizes the training procedure where it integrates the encoder-decoder architecture with the policy network in an end-to-end training framework.

IV. EXPERIMENTS

In our experiments we ask whether the proposed method can reliably handle inputs that lie within the data distribution yet were not observed during training, i.e., whether the policy can interpolate across realistic, unseen environment configurations produced by noisy measurements. In particular, we evaluate how well the framework mitigates distribution shift arising from disturbances in sensor or tracking readings. We benchmark against representative baselines on two simulated tasks (autonomous navigation of a ground vehicle and a Reach-Safe Franka Emika Panda manipulation task) and demonstrate results on Franka Panda hardware.

Our evaluation emphasizes safety metrics and robustness under randomized initial conditions, and in the following subsections we describe how we generate challenging test cases for thorough performance assessment.

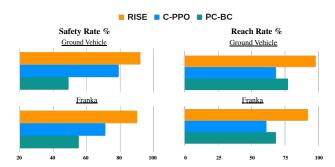


Fig. 3: Comparative study between all baselines and our approach based on evaluation matrics. The top plots illustrate the results for the Ground Vehicle Navigation setup, meanwhile, the bottom plots correspond to Franka Manipulator Task results from simulation.

A. Baselines and Evaluation Metrics

To evaluate our framework, we compare it against two baselines: Parameter-Conditioned Behavior Cloning (PC-BC) [13], which learns policies through behavior cloning with explicit conditioning on environmental parameters under randomized safety conditions, relying only on domain randomization for generalization. Primarily, we use the same inputs as our proposed approach, just without the Variational Encoder in this case; and Constrained Proximal Policy Optimization (C-PPO) [24], which extends PPO by incorporating safety constraints using Lagrange multipliers to penalize constraint violations during training. PC-BC tests whether including the variational encoder affects the performance or just a raw exposure to diverse conditions alone enables generalization, while C-PPO provides a reinforcement learningbased comparison that explicitly incorporates constraints. The performance of our method and its baselines are assessed using the following two key metrics:

Algorithm 1 Training Algorithm

```
Require: Dataset \mathcal{D} = \{(s, a, s', \text{obs})\}
  1: Stage 1: Pre-Train VAE (ELBO)
  2: for epoch = 1 to E_{\text{VAE}} do
  3:
             for batch \{obs\} from \mathcal{D} do
                   \mu, \sigma \leftarrow q_{\phi}(\text{obs})
  4:
                   \epsilon \sim \mathcal{N}(0, I), \ z = \mu + \sigma \odot \epsilon
  5:
                  \hat{\text{obs}} \leftarrow p_{\psi}(z)
  6:
                   \mathcal{L} \leftarrow -\log p_{\psi}(\hat{\text{obs}}|z) + \beta \text{ KL}[q_{\phi}(z|\text{obs})||p(z)]
  7:
                  Update \phi, \psi \leftarrow \nabla_{\phi, \psi} \mathcal{L}_{ELBO}
  8:
             end for
  9:
10: end for
11: Freeze VAE parameters: \phi, \psi \leftarrow detach
12: Stage 2: Train NN policy
      for epoch = 1 to E_{\pi} do
13:
             for batch \{(s, a, g, \text{obs})\} from \mathcal{D} do
14:
                  \mu, \sigma \leftarrow q_{\phi}(\text{obs})
                                                          ⊳ no gradients into VAE
15:
 16:
                  for m = 1 \dots M do \triangleright M unique perturbations
17:
                         \epsilon^{(m)} \sim \mathcal{N}(0, I), \ z^{(m)} = \mu + \sigma \odot \epsilon^{(m)}
18:
                        \hat{a}^{(m)} \leftarrow \pi_{\theta}(s, g, z^{(m)})
\mathcal{L}_{\pi} += \ell(\hat{a}^{(m)}, a)
19:
                                                                  ⊳ e.g., MSE / NLL
20:
                  end for
21:
                  \mathcal{L}_{\pi} \leftarrow \mathcal{L}_{\pi}/M
22:
                  Update \theta \leftarrow \nabla_{\theta} \mathcal{L}_{\pi}
23:
24.
25: end for
26: return trained policy \pi_{\theta} (VAE used in inference to
       sample z)
```

- 1) **Safety Rate**: Percentage of test trials in which the agent doesn't collide with the obstacle at any timestamp.
- 2) Reach Rate: Percentage of complete trials where the learned policy successfully reaches the goal. Note that a successfully reached episode is one during which the agent doesn't collide into the obstacle at any point in time during the entire episode.

B. Autonomous Navigation of a Ground Vehicle

In this task, the agent, an autonomous ground vehicle must reach the goal while navigating an environment containing a dynamic obstacle. The agent's state is represented as $s=(x,y,\theta)$, where (x,y) denotes position and θ is orientation. The action space consists of linear and angular velocities, $a=(v,\omega)$. The environment features a moving obstacle whose position is sampled to ensure a safe margin between the agent's initial state and goal. The obstacle's radius varies in a range of values to introduce variability, and its velocity is dynamically assigned to create unpredictable motion.

C. Franka Manipulator Task

In this task, a Franka Panda manipulator must reach a parameterized goal while avoiding obstacles in its workspace. For this experiment, we have used safe-panda-gym simulation environment [25], [26]. The action space comprises

end-effector displacement, a=(dx,dy,dz), applied through Position Control.

D. Training Data & Evaluation

Expert Data Generation: Expert demonstrations are generated using a mixture of experts like model predictive control [1], control barrier function (CBF), etc. The dataset, which includes 10k expert demonstrations, is constructed by randomly sampling initial robot states, and goal position.

Evaluation and Comparative Analysis: We evaluate performance across 1000 sampled test scenarios, all sampled by adding random noise (sampled from standard normal distribution) to the training data to emulate the desired noise in environmental parameter readings. Figure 3 summarizes the results. Our approach outperforms both C-PPO and PC-BC on both the evaluation metrics. It achieves the highest Safety Rate while maintaining a superior Reach Rate. Although C-PPO comes close in terms of safety, it struggles with goal-reaching performance, thus showing its conservative nature. On the other hand, we see PC-BC to be more aggressive and hence, suffers from frequent collisions. These results underscore the ability of RISE to balance safety while maintaining task performance.



Fig. 4: Illustration shows **Franka Panda manipulator** advancing toward its designated target (green region) while executing collision avoidance maneuvers in the presence of a dynamic obstacle (red sphere). Arrows indicate direction.

Hardware Results: We further validate our framework on a physical Franka Emika Panda manipulator. In this setup, virtual obstacles are employed, and environmental parameters (obstacle properties and goal locations) are provided in real time to the policy. Figure 4 presents key frames from the hardware demonstration, confirming successful goalreaching with effective obstacle avoidance.

V. CONCLUSION

In this paper, we proposed a practical imitation-learning framework that makes policies robust to realistic measurement uncertainty by conditioning them on a variational latent representation of environment parameters. By sampling plausible obstacle states from the VAE posterior during training, the policy learns to interpolate across nearby, realistic percepts and therefore behaves more conservatively and reliably in noisy, dynamic scenes without requiring explicitly requiring to train the model on these datapoints, nor requiring the exact dynamics models. Validation on an autonomous ground vehicle and a Franka Emika Panda demonstrates improved safety while preserving goal-reaching performance versus baselines.

REFERENCES

- [1] D. Mayne, J. Rawlings, C. Rao, and P. Scokaert, "Constrained model predictive control: Stability and optimality," *Automatica*, vol. 36, no. 6, pp. 789–814, 2000. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0005109899002149
- [2] S. Bansal, M. Chen, S. Herbert, and C. J. Tomlin, "Hamilton-jacobi reachability: A brief overview and recent advances," in 2017 IEEE 56th Annual Conference on Decision and Control (CDC). IEEE, 2017, pp. 2242–2253.
- [3] M. Tayal, A. Singh, S. Kolathaya, and S. Bansal, "A physics-informed machine learning framework for safe and optimal control of autonomous systems," 2025. [Online]. Available: https://arxiv.org/abs/2502.11057
- [4] A. D. Ames, X. Xu, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs for safety critical systems," *IEEE Transactions on Automatic Control*, vol. 62, no. 8, pp. 3861–3876, 2017.
- [5] M. Tayal, R. Singh, J. Keshavan, and S. Kolathaya, "Control barrier functions in dynamic uavs for kinematic obstacle avoidance: A collision cone approach," in 2024 American Control Conference (ACC). IEEE, 2024, pp. 3722–3727.
- [6] M. Tayal, A. Singh, P. Jagtap, and S. Kolathaya, "Semi-supervised safe visuomotor policy synthesis using barrier certificates," arXiv preprint arXiv:2409.12616, 2024.
- [7] J. Lee, C. Paduraru, D. J. Mankowitz, N. Heess, D. Precup, K.-E. Kim, and A. Guez, "Coptidice: Offline constrained reinforcement learning via stationary distribution correction estimation," in *International Conference on Learning Representations*, 2022.
- [8] M. Laskey, J. Lee, R. Fox, A. Dragan, and K. Goldberg, "Dart: Noise injection for robust imitation learning," in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., vol. 78. PMLR, 13–15 Nov 2017, pp. 143–156. [Online]. Available: https://proceedings.mlr.press/v78/laskey17a.html
- [9] Y. U. Ciftci, D. Chiu, Z. Feng, G. S. Sukhatme, and S. Bansal, "Safe-gil: Safety guided imitation learning for robotic systems," arXiv preprint arXiv:2404.05249, 2024.
- [10] A. Robey, H. Hu, L. Lindemann, H. Zhang, D. V. Dimarogonas, S. Tu, and N. Matni, "Learning control barrier functions from expert demonstrations," in 2020 59th IEEE Conference on Decision and Control (CDC), 2020, pp. 3717–3724.
- [11] M. Tayal, A. Singh, P. Jagtap, and S. Kolathaya, "Cp-ncbf: A conformal prediction-based approach to synthesize verified neural control barrier functions," arXiv preprint arXiv:2503.17395, 2025.
- [12] R. K. Cosner, Y. Yue, and A. D. Ames, "End-to-end imitation learning with safety guarantees using control barrier functions," in 2022 IEEE 61st Conference on Decision and Control (CDC). IEEE, 2022, pp. 5316–5322.
- [13] Y. Ding, C. Florensa, P. Abbeel, and M. Phielipp, "Goal-conditioned imitation learning," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2019/file/c8d3a760ebab631565f8509d84b3b3f1-Paper.pdf
- [14] M. Reuss, M. Li, X. Jia, and R. Lioutikov, "Goal-conditioned imitation learning using score-based diffusion policies," arXiv preprint arXiv:2304.02532, 2023.
- [15] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ser. ICML'17. JMLR.org, 2017, p. 22–31.
- [16] M. Alshiekh, R. Bloem, R. Ehlers, B. Könighofer, S. Niekum, and U. Topcu, "Safe reinforcement learning via shielding," *Proceedings* of the AAAI Conference on Artificial Intelligence, vol. 32, no. 1, Apr. 2018. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/ view/11797
- [17] D. A. Pomerleau, "Efficient training of artificial neural networks for autonomous navigation," *Neural computation*, vol. 3, no. 1, pp. 88–97, 1991.
- [18] S. Ross, G. J. Gordon, and J. A. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," 2011. [Online]. Available: https://arxiv.org/abs/1011.0686
- [19] A. Y. Ng and S. J. Russell, "Algorithms for inverse reinforcement learning," in *Proceedings of the Seventeenth International Conference*

- on Machine Learning. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2000, p. 663–670.
- [20] F. Castañeda, H. Nishimura, R. T. McAllister, K. Sreenath, and A. Gaidon, "In-distribution barrier functions: Self-supervised policy filters that avoid out-of-distribution states," in *Proceedings of The* 5th Annual Learning for Dynamics and Control Conference, ser. Proceedings of Machine Learning Research, N. Matni, M. Morari, and G. J. Pappas, Eds., vol. 211. PMLR, 15–16 Jun 2023, pp. 286–299. [Online]. Available: https://proceedings.mlr.press/v211/ castaneda23a.html
- [21] T. Schaul, D. Horgan, K. Gregor, and D. Silver, "Universal value function approximators," in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 1312–1320. [Online]. Available: https://proceedings.mlr.press/v37/schaul15.html
- [22] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," 2022. [Online]. Available: https://arxiv.org/abs/1312.6114
- [23] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," in *International* conference on learning representations, 2017.
- [24] A. Stooke, J. Achiam, and P. Abbeel, "Responsive safety in reinforcement learning by PID lagrangian methods," in *Proceedings* of the 37th International Conference on Machine Learning, ser. Proceedings of Machine Learning Research, H. D. III and A. Singh, Eds., vol. 119. PMLR, 13–18 Jul 2020, pp. 9133–9143. [Online]. Available: https://proceedings.mlr.press/v119/stooke20a.html
- [25] S. W. Tosin Oseni, "Safe panda gym," https://github.com/tohsin/ Safe-panda-gym, 2022.
- [26] Q. Gallouédec, N. Cazin, E. Dellandréa, and L. Chen, "panda-gym: Open-Source Goal-Conditioned Environments for Robotic Learning," 4th Robot Learning Workshop: Self-Supervised and Lifelong Learning at NeurIPS, 2021.