

Continual Domain Randomization

Josip Josifovski^{1*}, Sayantan Auddy^{2*}, Mohammadhossein Malmir¹,
Justus Piater^{2,4}, Alois Knoll¹ and Nicolás Navarro-Guerrero³

Abstract—Domain Randomization (DR) is commonly used for sim2real transfer of reinforcement learning (RL) policies in robotics. Most DR approaches require a simulator with a fixed set of tunable parameters from the start of the training, from which the parameters are randomized simultaneously to train a robust model for use in the real world. However, the combined randomization of many parameters increases the task difficulty and might result in sub-optimal policies. To address this problem and to provide a more flexible training process, we propose *Continual Domain Randomization* (CDR) for RL that combines domain randomization with continual learning to enable sequential training in simulation on a subset of randomization parameters at a time. Starting from a model trained in a non-randomized simulation where the task is easier to solve, the model is trained on a sequence of randomizations, and continual learning is employed to remember the effects of previous randomizations. Our robotic reaching and grasping tasks experiments show that the model trained in this fashion learns effectively in simulation and performs robustly on the real robot while matching or outperforming baselines that employ combined randomization or sequential randomization without continual learning. Our code and videos are available at <https://continual-dr.github.io/>.

Index Terms—Domain randomization, sim2real transfer, continual reinforcement learning, robotic manipulation

I. INTRODUCTION

Deep reinforcement learning (DRL) in robotics research relies heavily on simulated environments, mainly because of the sample inefficiency of DRL algorithms and the absence of safe exploration guarantees, which make it impractical to train such algorithms on physical robots directly. While simulations provide virtually unlimited training data and eliminate safety concerns during training, their approximative nature can significantly degrade the model’s performance on the real system in sim2real transfer [1] – a phenomenon known as the *reality gap*.

To bridge the reality gap, one can try to attain *zero-shot* sim2real transfer by crafting an accurate simulation and perfectly mimicking the behavior of the real system [1]. Unfortunately, even for elementary systems, simulating reality with sufficient precision is rarely possible. Alternatively, a simulation-trained DRL model can be *finetuned* on the target to improve performance. This *domain adaptation* approach updates the policy learned in simulation and usually leads to

improved performance. However, it makes the policy system-specific and requires at least some real-world training, which can be unsafe and time-consuming [2].

Due to the impracticality of developing perfect simulations and the limitations of domain adaptation, recent works on *domain randomization* [3] rely on the randomization of simulation parameters. Here, instead of accurately simulating every parameter of the real system, the relevant simulation parameters are randomized to cover the difference between the real world and simulation. However, it can be challenging to precisely identify all the essential randomization parameters at the start. Moreover, the excessive randomization of many parameters jointly increases the task complexity and the agent’s uncertainty, making it difficult to find a good policy in simulation [4].

To address these limitations, we propose *Continual Domain Randomization* (CDR), which combines *domain randomization* [3] with *continual learning* [5] (CL) for sim2real transfer. Viewing domain randomization as a CL problem provides greater flexibility in the initial choice of the randomization parameters. Additionally, randomizing a subset of all parameters at a time reduces the risks of excessive randomization. In CDR, we start with an agent trained in an idealized simulation (with all randomizations disabled) and successively train the agent in differently randomized simulations. We treat each randomization as a separate CL *task* and enable the DRL agent to continually train under each randomization without *catastrophically forgetting* [5] the knowledge gained from previous training runs (i.e., past *tasks*). We implement an offline and an online version of CDR by combining the Proximal Policy Optimization (PPO) [6] algorithm with the regularization-based CL algorithm *Elastic Weight Consolidation* (EWC) [7] and its online variant [8]. Nevertheless, CDR is a general strategy that can be implemented with other DRL or CL algorithms.

We perform experiments for sim2real transfer of two different robotics tasks, reaching and grasping, and compare the real-world performance of policies learned using CDR against *finetuning*, *full-randomization*, and the ideal model trained without randomization. Our results show that CDR can match or outperform baselines that jointly randomize all parameters throughout the training process. At the same time, it is more stable and less sensitive to the order of randomizations than sequential baselines like *finetuning*, which only optimizes the policy for the newest randomization and ignores previous ones. These results highlight the potential of CDR as a general and flexible framework for helping bridge the reality gap in DRL-based robotics.

¹ School of Computation Information and Technology, Technical University of Munich, Germany.

² Department of Computer Science, University of Innsbruck, Austria.

³ L3S Research Center, Leibniz Universität Hannover, Germany.

⁴ Digital Science Center (DiSC), University of Innsbruck, Austria.

* Equal contribution

Sayantan Auddy is supported by a doctoral scholarship granted by the University of Innsbruck, Vice-Rectorate for Research.

II. RELATED WORK

Even though there are some cases of sim2real transfer where measuring and precisely simulating all relevant parameters of the real system may be possible [9], [10], due to wear and tear in the robot’s hardware or changes in external conditions like temperature or humidity, they can vary significantly [1]. Intuitively, providing some degree of variability in the simulator parameters and training a model that can deal with a range of parameter values would lead to a more robust sim2real transfer [11]. For this reason, *domain randomization* (DR) [3], [12]–[18] has become a widely used sim2real technique. However, inappropriate and excessive randomization can increase the uncertainty of the system and the difficulty of solving the task in simulation, leading to sub-optimal policies and longer training times [4].

One approach to circumvent those problems is *automatic domain randomization* [14]. Here, the randomization increases progressively as the policy under the current randomization improves the performance. In contrast, *active domain randomization* [15] starts from wide randomization ranges and learns the most informative ranges for training a robust policy. While the previous approaches assume that only simulated data is available, data from the real system can be used in offline fashion [17], [19] to find suitable parameter distributions, or to identify which parameters are relevant for sim2real transfer through causal discovery [18].

Continual learning (CL) algorithms enable knowledge accumulation from a sequence of tasks without *catastrophic forgetting* [5]. CL has traditionally targeted supervised learning scenarios using static datasets or reinforcement learning problems with discrete actions [7], [8]. More recently, CL has been applied to more challenging scenarios such as robot learning, both in simulation [20] as well as in the real world [21], [22]. Continual reinforcement learning (CRL) in continuous action and observation spaces [23] has also emerged as a popular means of training robots to learn multiple tasks. CRL scenarios can be of different types [23]: (i) the state transition and action space for different tasks remain the same, but the reward function changes (the agent needs to perform fundamentally different actions such as pushing, reaching, etc., for each task) [24]; (ii) the state transition function changes, but the reward function and action space remain the same. Here, the environment dynamics change, but the ultimate goal of the agent remains the same [25]. The sim2real scenario belongs to the latter.

Progressive networks [26] have also been used for continual sim2real transfer [27]. For instance, as shown by Rusu et al., [27], an initial neural network can be trained in simulation to solve a task, and a new network, with lateral connections to the previous network for feature reusing, is instantiated and trained on the real robot. Another approach used for sim2real transfer is *Policy distillation* [28]. For instance, Traoré et al., [29] trained a separate policy and generated a dataset of observation-action pairs for each task in simulation. The datasets are then used in a supervised learning fashion to train a single model that can solve the tasks in the real environment.

Unlike other domain randomization approaches, in CDR no fixed set of randomization parameters during training is assumed – the model can train on a subset of randomization parameters at a time and adapt to new randomization parameters later. This is achieved by combining sequential randomization with continual learning, which retains the effects of previous randomizations and avoids overfitting to the last set of parameters. Unlike other continual sim2real approaches, CDR uses a single neural network trained across all tasks. Additionally, the online CDR variant avoids the linear increase in memory for keeping a separate dataset or network parameters per task. Finally, CDR can easily be combined with methods like *automated* [14] or *active* [15] DR to find suitable ranges for the subset of randomization parameters within a single training run in the sequence.

III. METHODOLOGY

A. Problem Description

Let $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, S_0, \gamma, T)$ define a Markov Decision Process (MDP) where \mathcal{S} denotes the state space, \mathcal{A} is the action space, $P(s_{t+1} \in \mathcal{S} | s_t \in \mathcal{S}, a_t \in \mathcal{A})$ is the state transition probability, R is the reward function, S_0 defines an initial state distribution, γ is the reward discount factor and T defines the finite horizon. The objective of reinforcement learning is to find the parameters θ for a parameterized policy $\pi_\theta(a_t | s_t)$ to maximize the expected return $\mathbb{E}_{\pi_\theta}[\sum_{t=0}^T \gamma^t R(s_t, a_t)]$ where $s_0 \sim S_0$, $s_{t+1} \sim P$ and $a_t \sim \pi_\theta(a_t | s_t)$. Let Φ denote a real system with unknown dynamics, and $\mathcal{E} = \{\epsilon_1, \dots, \epsilon_n\}$ denote the set of all randomization parameters of a simulator Σ of Φ . The goal is to optimize θ using Σ such that it maximizes the expected return on Φ (zero-shot sim2real transfer).

B. Continual Domain Randomization

In our suggested approach, *Continual Domain Randomization* (CDR), the policy parameters $\theta = \tilde{\theta}$ are initialized and optimized in an idealized simulation Ψ_\emptyset (with all randomizations disabled). Afterward, we sequentially optimize the resulting θ with CL on a sequence of simulations where a subset of randomization parameters is enabled each time, as shown in Fig. 1: $\tilde{\theta} \xrightarrow{\Psi_\emptyset} \theta_\emptyset \xrightarrow{\Psi_{\epsilon_1}} \theta_1 \dots \xrightarrow{\Psi_{\epsilon_n}} \theta_n$.

CDR- λ : combines the PPO [6] algorithm with the popular regularization-based CL algorithm Elastic Weight Consolidation (EWC) [7], λ denotes the regularization process with EWC. EWC is an approximate Bayesian method that quantifies the importance of neural network parameters for remembering previous knowledge using the empirical *Fisher information matrix* (FIM). The FIM is used to estimate the posterior distribution over the network parameters, approximating it as a multivariate Gaussian centered around the maximum a posteriori (MAP) estimate. The diagonal of the covariance matrix is then used to determine parameter importance. The loss function for learning task C after tasks A and B is:

$$\mathcal{L}(\theta) = \mathcal{L}_C(\theta) + \frac{\lambda_B}{2} \|\theta - \theta_B^*\|_{F_B}^2 + \frac{\lambda_A}{2} \|\theta - \theta_A^*\|_{F_A}^2 \quad (1)$$

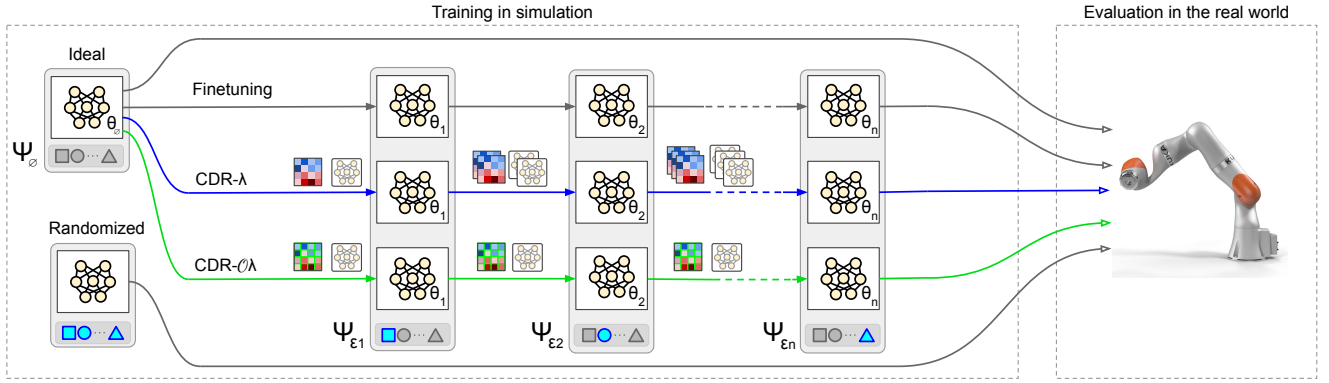


Fig. 1. Overview of our proposed CDR approach. CDR- λ (blue arrow) uses a set of network snapshots and Fisher matrices (one for each past task) for continual learning, while CDR- $\mathcal{O}\lambda$ (green arrow) uses only a single parameter snapshot and Fisher matrix. Other baselines are shown with gray arrows. Each shape ($\square, \circ, \dots, \triangle$) represents a unique randomization parameter set. Disabled and enabled sets are indicated with gray and blue respectively.

where \mathcal{L}_C is the task-specific loss for task C , λ_A and λ_B are the regularization constants, θ_A^* and θ_B^* are the saved MAP parameters after learning tasks A and task B respectively, and F_A and F_B are the diagonal Fisher information matrices for the first two tasks. In this paper, we adapt PPO for CL by regularizing the parameters of *only* the actor network with EWC. After learning each task, we roll out trajectories of the PPO agent trained on the most recent environment and store these in a replay buffer. After that, we compute the (diagonal) empirical Fisher information for the latest task using the squared gradients of the log-likelihoods of the action [30] predicted by the PPO policy:

$$F(\theta) = \frac{1}{n} \sum_i^n \nabla_{\theta} \log p_{\theta}(y_i|x_i) \nabla_{\theta} \log p_{\theta}(y_i|x_i)^T \quad (2)$$

here, n denotes the number of samples in the replay buffer, x_i and y_i denote the input and output of the policy network, and θ denotes the parameters of the policy network. Since the critic does not contribute to the log-likelihood computation, the parameters of the critic are not regularized. In contrast to previous implementations of PPO and EWC [31], but similar to online-EWC [8], we normalize the elements of F to be in the range $[0.0, 1.0]$ after each task. This normalization decouples the regularization constant λ from the arbitrary scaling of the Fisher information of individual tasks and allows us to use of a single value of λ for all tasks.

A disadvantage of regular EWC is that for each randomization (task), a separate snapshot of the actor network parameters and the corresponding empirical diagonal Fisher information needs to be saved. This results in the undesired linear scaling of storage requirements. To counter this, we suggest combining PPO with online-EWC [8] creating **CDR- $\mathcal{O}\lambda$** . After learning a task, a replay buffer is filled, and the empirical Fisher information is computed with Eq. 2. Thus, instead of maintaining a separate copy of the network and Fisher information for each task, only one copy of the network after the most recent task is kept, and the network is updated using

$$\mathcal{L}_i(\theta) = \tilde{\mathcal{L}}_i(\theta) + \frac{\lambda}{2} \|\theta - \theta_{i-1}^*\|_{\gamma F_{i-1}^*}^2 \quad (3)$$

where $\mathcal{L}_i(\theta)$ is the overall loss for the current task, $\tilde{\mathcal{L}}_i(\theta)$ is the task-specific PPO loss, θ_{i-1}^* is the snapshot of the network after the previous task, and F_{i-1}^* is the *online* Fisher information after the previous task. F_{i-1}^* and the Fisher information of the current task F_i are used to update the online Fisher information F_i^* using a hyperparameter γ :

$$F_i^* = \gamma F_{i-1}^* + F_i \quad (4)$$

We again normalize F_i to be in the range $[0.0, 1.0]$ to avoid any arbitrary scaling. Similar to CDR- λ , we only regularize the PPO actor network for CDR- $\mathcal{O}\lambda$.

C. Baselines

We implement multiple baselines [4] with PPO [6] (Fig. 1):

- (i) **Ideal**: Training is conducted in the *ideal* simulation environment without any randomization.
- (ii) **Finetuning**: The model trained in the ideal simulation is used as a starting point and then successively *finetuned* on each of the randomization parameters, one at a time, without any continual learning.
- (iii) **Randomized**: The model is trained by activating all the randomization parameters together at the same time.

IV. EXPERIMENT SETUP

For the experiments, we consider the *reach-and-balance* [4] task with a 2-joint-controlled robotic arm, and a robotic *grasping* task with a kinematically redundant, 7-joint-controlled arm and a three-fingered gripper [32]. Both tasks are defined in continuous state and action spaces, and the RL agent exercises fine-grained direct control over the robot's joint velocities. We intentionally use joint space control, because it requires minimal information about the real system's dynamics [33]. Moreover, in this case the agents are supposed to learn to compensate for dynamic effects, e.g., inertial forces [33] through randomization, which is particularly important and part of our sim2real transfer analysis.

A. Reacher Task

The *reacher* task is defined as:

$$\mathbf{s}_t = [\Delta x_t, \Delta y_t, \Delta z_t, q_1, q_2] \in \mathbb{R}^5 \quad (5)$$

$$\mathbf{a}_t = [\dot{q}'_1, \dot{q}'_2] \in \mathbb{R}^2 \quad (6)$$

$$r_t = \begin{cases} -d_t & \text{if } C \text{ is false} \\ -d_t \times (T - t) & \text{otherwise} \end{cases} \in \mathbb{R} \quad (7)$$

where \mathbf{s}_t is the agent's state at timestep t , Δx , Δy and Δz are the distances between the end-effector and the target along the corresponding coordinate axes, and q_i are the joint positions. \mathbf{a}_t is the set of joint velocities the agent commands at timestep t . The reward r_t at each timestep is defined as the negative Euclidean distance d_t between the end-effector and the target, unless the robot collides with the floor or reaches joint limits, in which case the terminal condition C becomes true, the episode terminates and the reward is the negative distance at that timestep multiplied with the remaining timesteps to the maximal episode length T .

B. Grasper Task

The *grasping* task is defined as:

$$\mathbf{s}_t = [\Delta x_t, \Delta y_t, \Delta z_t, q_1, q_2, q_3, q_4, q_5, q_6, q_7, \Delta f x_t, \Delta f y_t, \Delta f z_t, \Delta u x_t, \Delta u y_t, \Delta u z_t] \in \mathbb{R}^{16} \quad (8)$$

$$\mathbf{a}_t = [\dot{q}'_1, \dot{q}'_2, \dot{q}'_3, \dot{q}'_4, \dot{q}'_5, \dot{q}'_6, \dot{q}'_7] \in \mathbb{R}^7 \quad (9)$$

$$r_t = \begin{cases} -dp_t - 0.5 \times du_t - c & \text{if } t < T \\ -dp_t - 0.5 \times du_t - c + g & \text{if } t = T \end{cases} \in \mathbb{R} \quad (10)$$

where \mathbf{s}_t is the state of the agent at timestep t , Δx_p , Δy_p and Δz_p are the distances between the end-effector and the center of the object along the corresponding coordinate axes, and q_{1-7} are the joint positions. $\Delta f x$, $\Delta f y$, $\Delta f z$ are the distances between the *forward* unit direction vector of the object and the *forward* unit direction vector of the end-effector along each axis, and $\Delta u x$, $\Delta u y$, $\Delta u z$ are the distances between the *up* unit direction vector of the plane on which the object is placed and the *up* unit direction vector of the end-effector along each axis. \mathbf{a}_t is the set of joint velocities the agent commands at time t . Each episode has a fixed duration T . If a collision occurs or the robot reaches joint limits during the episode, zero velocities are commanded until the episode ends. At the end of the episode, if there were no collisions and the joint limits were not reached, the robot stops moving and the gripper closes in an attempt to grasp the object. The reward at timestep t is defined in terms of the negative Euclidean distance dp_t between the end-effector and the center of the object, the negative Euclidean distance du_t between the *up* unit vectors of the gripper and the floor, and a penalty coefficient c that is 0 by default and 1 if there is a collision or joint limits are reached during the episode. At the end of the episode, if the grasping attempt was successful, the term g is set to 10, otherwise it is 0.

C. Robot Platform

The simulated and real-world environments are shown in Fig. 2. We use the Unity-based platform introduced in [25] for the simulation. For both *reacher* and *grasper*, we use the *KUKA LBR iiwa 14* robotic manipulator [34], and for the grasping task we use the 3-Finger Adaptive Robotic Gripper

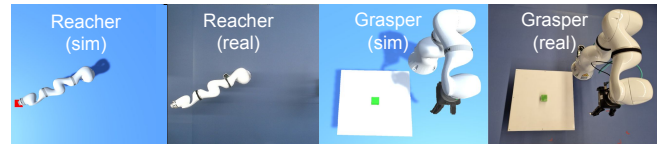


Fig. 2. The simulated and real environments for reaching and grasping.

(3F model) by Robotiq in pinch mode, relying on the gripper parameters and meshes from the URDF data provided by the ROS-Industrial [35] package. The real robot is controlled with the Iiwa stack [36] via Robot Operating System (ROS) [37]. The box used for the grasping experiments is a cube with a side of 50mm.

D. Training Procedure

We use three randomisation parameters commonly used in the literature [4], [13]: (i) Latency L : random delay in observing the current state; (ii) Torque T : the joint stiffness and damping coefficients are randomized, resulting in a randomized torque for reaching the commanded velocity; and (iii) Noise N : noise drawn from a uniform distribution is randomly incorporated to simulate sensor inaccuracies. The implementation is identical to [4] and the ranges of the randomization parameters are provided in Tab. I.

The efficacy of sim2real transfer under sequential training strategies like finetuning or CDR may be influenced by the ordering of randomization parameters, as each parameter can have varying effects on the transfer process [4]. Instead of unnecessarily testing all possible permutations of L , N , and T , we empirically evaluate the importance of each parameter for sim2real transfer on the reaching task to determine the best and worst case ordering scenario. In the best scenario, the parameter that leads to the best sim2real transfer is at the end of the training sequence, while in the worst case scenario, it is at the beginning. For the empirical evaluation, we perform training on a single-parameter randomized simulation from scratch and evaluate the models in ideal simulation and on the real system. The results from 5 independent seeds for each parameter, as well as the baselines where no parameter or all parameters are randomized, are presented in Fig. 3. Based on the results, we define *TLN* and *NLT* orderings for the sequential training strategies in the experiments.

We extend the PPO [6] implementation in *StableBaselines3* [38] to implement the methods described in Sec. III.

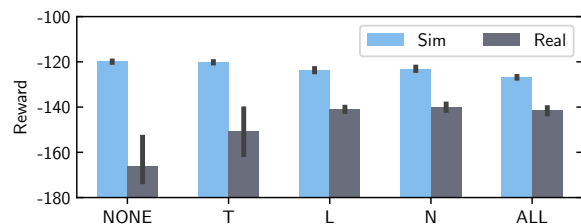


Fig. 3. Effects of different randomization parameters on sim2real transfer for the reaching task.

TABLE I
TRAINING HYPERPARAMETERS.

Parameter (PPO)	Value	Parameter	Value
Bias-var tradeoff (λ)	0.95	Torque (T)	stiff: $[10^1, 10^3]$ damp: $[10^1, 10^3]$
Disc. factor (γ)	0.99	Latency (L)	$[0, 1]$ sec
Learning rate (α)	25×10^{-5}	Noise (N)	$[0, 10]$ %
Val. fun. coeff. (c_1)	0.5	EWC reg. const. λ	5×10^3
Ent. coeff. (c_2)	0.0	EWC buffer limit	100
Clip range (ϵ)	0.1	EWC replay batch	32
Epochs	10	EWC replay samples	5000
Mini-batches	32	online-EWC γ	0.95

We use the default network architecture (2-layer MLP, 64 units/layer, *tanh* activations). Due to the different state and action space dimensions, the network for *grasper* is much larger than that for *reacher*. We use the same PPO-specific hyperparameters as in [4]. For the EWC hyperparameters, we start from literature-recommended values [7] and adjust them empirically. All hyperparameters are listed in Tab. I.

For *reacher*, we use an episode length of 500 timesteps, with a timestep duration of 20ms (50Hz control frequency), for a total simulated time of 10s per episode. The maximal allowed joint speed is 20 deg/sec. For *grasper*, due to safety concerns on the real system, we reduce the maximum joint speed to 10 deg/sec, adapt the timestep duration to 50ms (20Hz control frequency), and the episode length to 200 timesteps, for the same total simulated time per episode of 10s. Under the *ideal* training strategy, the simulator has all randomization parameters disabled throughout the entire training duration, i.e., there is no latency or noise, and the commanded joint velocities are reached instantaneously. Under the fully randomized strategy, all randomization parameters are simultaneously enabled. We start with a model pre-trained in ideal simulation for the CDR and *finetuning* strategies. Then, we train them sequentially on each of the N , L , and T -only randomized simulations for equal duration. For the reaching task, the total training duration is 4×10^6 timesteps, while for the grasping task it is 10×10^6 timesteps due to the higher task complexity and a larger number of trainable parameters.

E. Evaluation Procedure and Metrics

While the training takes place only in simulation, evaluation takes place both in simulation (ideal setting) and on the real robot, see Fig. 2. In both cases, the same episode duration, control frequency, robot’s starting configuration, maximal joint speeds, joint limits, and object/target positions per task are used, leading to identical conditions in the simulation and the real system, such that any difference in performance of the evaluated model can be attributed to the sim2real gap. We report multiple metrics detailed in Tab. II, both in simulation and the real system.

TABLE II
EVALUATION METRICS.

Cumulative episodic reward (r_{ep}): The sum of the rewards collected by the agent during an evaluation rollout of T steps.

$$r_{\text{ep}} = \sum_{t=0}^T r_t \quad (11)$$

Continuity cost (\mathcal{C}) [39]: A measure of the smoothness of motion.

$$\mathcal{C} = 100 \times \mathbb{E}_t \left[\left(\frac{a_{t+1} - a_t}{\Delta_{\text{max}}^a} \right)^2 \right] \quad (12)$$

where a_{t+1} and a_t are consecutive actions, and Δ_{max}^a denotes the maximum difference between consecutive actions during the evaluation rollout. \mathcal{C} ranges between 0 (constant actions) and 100 (each consecutive action reaches opposing limits in the action space). Lower values denote a smoother robot motion and are desirable in the real system.

Distance to the target (d_{tgt}): The average Euclidean distance of the center of the robot’s end-effector p_t to the target p_{tgt} in the second half of the episode, when it is expected that the target is already reached. This is related to its stabilization strength [40], i.e., the ability to balance and avoid jittering motions around the target once it is reached. Lower values denote better performance.

$$d_{\text{tgt}} = \frac{2}{T} \sum_{t=T/2}^{t=T} \|p_t - p_{\text{tgt}}\|^2 \quad (13)$$

Grasping success (\mathcal{G}): Measures the percentage of successful grasps during multiple evaluation rollouts. Values range from 0 to 1, where one means that every grasp attempt was successful.

V. EXPERIMENT RESULTS

A. Reacher

The training progress of agents trained under different strategies and randomization orders to solve the *reaching* task is presented in Fig. 4. Each model is trained ten times with independent seeds and their performance throughout training is evaluated in non-randomized simulation. The *reaching* task is relatively easy to solve, especially in the absence of randomizations, which is why it takes less time for the *Ideal* baseline to converge w.r.t. the *Randomized* one. For the sequentially trained models the ordering of the randomization parameters greatly affects the extent to which the policy can be improved. The policy improves faster when torque randomization is first (TLN) due to the lower uncertainty in the agent’s state compared to the sequence starting with noise randomization (NLT). *Finetuning* adapts much faster than the CL models ($CDR-\lambda$, $CDR-\mathcal{O}\lambda$) and reaches performance in simulation similar to the *Ideal* model. Among the CL models, $CDR-\lambda$ changes more slowly and reaches a lower performance in simulation compared to $CDR-\mathcal{O}\lambda$ since there is a separate regularization for each randomization parameter under EWC training, which reduces the ability of the policy to change. $CDR-\mathcal{O}\lambda$ adapts faster and converges to a similar performance as the fully-randomized models.

The evaluation results in a non-randomized simulation and on the real system for each randomization strategy for the *reaching* task at the end of training are presented in Tab. III. The *Ideal* model performs best in simulation. However, it overfits to the simulation, and this leads to poor transfer

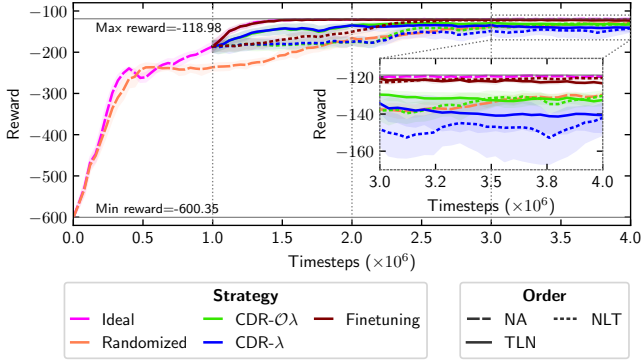


Fig. 4. Training progress for the reaching task. *Finetuning*, CDR- λ and CDR- $\mathcal{O}\lambda$ start from the *Ideal* model at 10^6 timesteps, and are sequentially trained on each randomization for 10^6 steps, shown by vertical dotted lines. The max reward is the best reward achieved by any agent, and the min reward is the one achieved by an agent that executes random actions.

TABLE III

REACHER EVALUATION IN SIMULATION AND REALITY. METRICS ARE EPISODIC REWARD r_{ep} , DIST. TO TARGET d_{tgt} , AND CONTINUITY COST \mathcal{C} .

Strategy	$r_{ep} \uparrow$		$d_{tgt} \downarrow$		$\mathcal{C} \downarrow$	
	Sim	Real	Sim	Real	Sim	Real
CDR- λ (NLT)	-140.59 \pm 15.21	-157.24 \pm 17.23	0.06	0.07	0.00	0.00
CDR- λ (TLN)	-140.37 \pm 15.15	-157.72 \pm 18.59	0.06	0.07	0.00	0.00
CDR- $\mathcal{O}\lambda$ (NLT)	-128.82 \pm 4.88	-143.88\pm4.82	0.02	0.02	0.00	0.00
CDR- $\mathcal{O}\lambda$ (TLN)	-130.97 \pm 4.26	-146.95 \pm 4.18	0.03	0.03	0.00	0.00
Finetuning (NLT)	-119.57 \pm 0.29	-197.78 \pm 7.23	0.00	0.19	0.01	0.16
Finetuning (TLN)	-121.86 \pm 1.09	-196.75 \pm 8.24	0.01	0.18	0.00	0.11
Ideal	-119.27\pm0.38	-208.46 \pm 2.17	0.00	0.20	0.01	0.30
Randomized	-129.16 \pm 1.94	-146.52 \pm 3.11	0.02	0.02	0.00	0.00

to the real system, resulting in lower rewards (r_{ep}), high continuity cost (\mathcal{C}), and unstable jittering motions around the target which results in high average distance to the target (d_{tgt}). *Randomized* has lower rewards in simulation but achieves very good sim2real transfer, leading to higher rewards, low continuity cost, and lowest distance to target on the real system. The CDR models also achieve very good sim2real transfer and perform comparably to or better than the fully *Randomized* baseline. Due to the stronger EWC regularization during training, CDR- λ reaches lower rewards and higher distance to target w.r.t. the CDR- $\mathcal{O}\lambda$ strategy, while it has comparable continuity cost to CDR- $\mathcal{O}\lambda$. Overall, CDR- $\mathcal{O}\lambda$ achieves the best rewards on the real system and has distance and continuity cost values similar to the fully *Randomized* baseline. *Finetuning* performs comparably to *Ideal* in simulation but better on the real system. However, due to the lack of any regularization, it finds over-confident policies in simulation, leading to much worse sim2real transfer than the CDR strategies. The qualitative differences in performance between the different strategies are also noticeable in the supplementary video.

B. Grasper

The results of the training progress of the different strategies for the *grasping* task evaluated in the non-randomized simulation are presented in Fig. 5. Each experiment is repeated five times with independent seeds. Compared to *reaching*, the *grasping* task is more challenging to solve due to the larger model sizes, the interaction with the object, and the

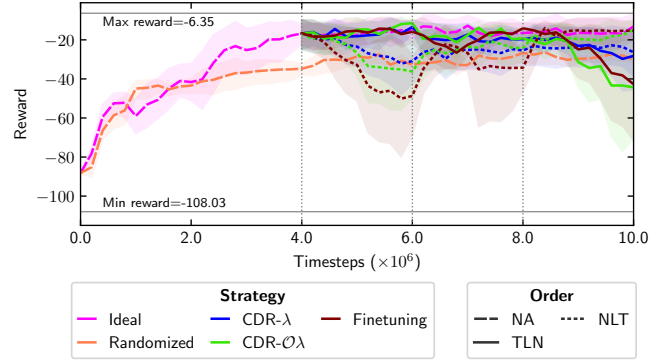


Fig. 5. Training progress for the grasping task. *Finetuning*, CDR- λ and CDR- $\mathcal{O}\lambda$ start from the *Ideal* model at 4×10^6 timesteps, and are sequentially trained on each randomization for 2×10^6 steps (vertical dotted lines). The max reward is the best reward achieved by any agent, and the min reward is the one achieved by an agent that executes random actions.

TABLE IV

GRASPER EVALUATION IN SIMULATION AND REALITY. METRICS ARE EPISODIC REWARD r_{ep} , CONTINUITY COST \mathcal{C} , AND GRASPING SUCCESS \mathcal{G} .

Strategy	$r_{ep} \uparrow$		$\mathcal{C} \downarrow$		$\mathcal{G} \uparrow$	
	Sim	Real	Sim	Real	Sim	Real
CDR- λ (NLT)	-25.84 \pm 15.10	-24.27 \pm 7.96	0.06	0.15	0.40	0.20
CDR- λ (TLN)	-27.63 \pm 20.25	-23.25 \pm 6.89	0.08	0.11	0.40	0.20
CDR- $\mathcal{O}\lambda$ (NLT)	-12.74 \pm 6.34	-22.78 \pm 21.97	0.05	0.12	0.78	0.60
CDR- $\mathcal{O}\lambda$ (TLN)	-44.81 \pm 32.47	-23.61 \pm 6.11	0.08	0.09	0.23	0.20
Finetuning (NLT)	-12.96 \pm 5.99	-14.99\pm5.56	0.07	0.09	0.59	0.40
Finetuning (TLN)	-42.88 \pm 35.37	-27.24 \pm 15.30	0.08	0.02	0.33	0.40
Ideal	-12.43\pm7.69	-32.90 \pm 41.40	0.06	0.45	0.77	0.60
Randomized	-28.42 \pm 6.00	-25.89 \pm 7.28	0.04	0.01	0.20	0.40

high precision needed for successful grasping.

With the *Ideal* model, grasping the object and achieving a high reward is easier, while with the fully *Randomized* model, it is more difficult due to the high uncertainty of the agent's state. Among the sequential training strategies, it can be seen that the task is significantly more challenging to solve when the noise randomization N is encountered in the *TLN* or *NLT* sequences. The policy found at the previous randomization step significantly changes to adapt to the noise-randomized simulation. Overall, *Finetuning* has higher deviations during training, while the deviations are lower and more stable throughout training for the CDR strategies, with CDR- λ being more consistent in comparison to CDR- $\mathcal{O}\lambda$. This is likely due to the separate regularization for each randomization parameter.

The evaluation results for *grasping* in ideal simulation and the real world after training in the simulation are presented in Tab. IV. Similar to the *reaching* task, *Ideal* achieves high reward in simulation but achieves low average reward and high continuity cost (unsafe jittering behavior, see supplementary video) on the real system, even though the object is grasped in most runs. *Randomized* has moderate success in finding a good policy in simulation and has the lowest grasp success in non-randomized simulation. However, it has good sim2real transfer and achieves the lowest continuity cost. From the sequential strategies, we notice that they are sensitive to the ordering and achieve higher success rates but also higher continuity cost when torque is randomized last in the sequence (*NLT*), and lower continuity cost but also lower grasp

success when noise is at the end of the sequence (*TLN*). If we look at the rewards on the real system, *Finetuning* deviates significantly and, depending on the sequence, can achieve the best or worst reward from any model that involves randomization, with a difference of about 12%. On the other hand, the CDR models are less susceptible to the ordering of randomization, as they achieve similar rewards regardless of the sequence. Moreover, the rewards of the CDR models lie within the range of the best and worst *Finetuning* strategy and are better than the reward of the *Randomized* strategy.

C. Effect of Continual Learning Regularization Strength

In the context of regularization-based CL methods used to implement CDR, the regularization strength λ plays a key role. A very low value makes the CDR models behave like *Finetuning*, while a very high value makes it difficult to adapt to new randomizations after being initially trained in ideal simulation. Taking this into account, in our main experiments we have used a value for λ solely based on empirical evaluation in simulation in order to maintain the zero-shot sim2real condition. Here, we perform additional experiments to provide better intuition on a suitable range for λ .

For this study, we employ the same training procedure as in the *reaching* experiment in Sec. V-A. We use the same hyperparameters (Tab. I) and only change the value of the EWC λ parameter, testing a range of different values of $\lambda = x \times 10^y$ with $x \in \{1, 5\}$ and $y \in \{0, \dots, 4\}$. We train models of *CDR- λ* and *CDR- $\mathcal{O}\lambda$* in the ideal simulation followed by two different randomization sequences *TLN* and *NLT*. We repeated them three times with independent seeds. The evaluation results of the models in the ideal simulation and on the real system are shown in Fig. 6.

CDR models perform well in simulation for lower values of λ . However, they become sensitive to the order of randomizations and transfer poorly to the real system. On the other hand, higher values of λ prevent changes to the model pre-trained on non-randomized simulation used as a starting point and perform poorly both in simulation and on the real system. In this case, the continuity cost on the real system is low because the agent performs slow moves toward the target but does not learn to balance and overshoots the target. The middle ranges of λ perform reasonably well in simulation and give the best results on the real system while being relatively robust to the order of randomizations.

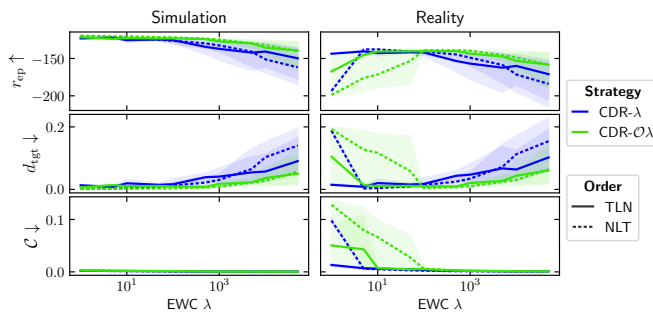


Fig. 6. Effect of the EWC regularization constant λ on episodic reward (r_{ep}), distance to target (d_{tgt}) and continuity cost (C) for sim2real transfer.

In general, full randomization of many parameters simultaneously makes the task more difficult [4], [14], [15], and our experiments show that it takes longer for the models to converge under this setting. Sequential randomization starting from a model pre-trained on non-randomized simulation offers a middle ground between no randomization and full randomization. However, the ordering of the randomizations can affect sim2real transfer. Unfortunately, in zero-shot transfer (without testing on the target system), the effect of each randomization parameter cannot be known with certainty. This can be detrimental since common sequential randomization strategies overfit to the latest randomization they were trained on.

CDR is less susceptible to the order of randomizations and can transfer the effects of earlier randomizations in the sequence to the final model. This can be done either explicitly with *CDR- λ* where the transfer comes at the price of linear increase of the memory, or implicitly with *CDR- $\mathcal{O}\lambda$* where memory does not increase but the regularization effect for continual learning is weaker. This paper uses CDR for *zero-shot* sim2real transfer. However, CDR can also be adapted to provide more flexible *domain adaptation* behavior by viewing the real system as the final randomization in the training sequence. EWC- λ can be set to 0 so that the policy is finetuned to the real system dynamics or to values greater than 0 in cases where the system parameters can vary over time due to wear and tear or external conditions and the policy needs to stay more general.

A limitation of our current training approach on one randomization at a time is that we factorize the space of randomization parameters. The complex interactions between separate randomizations in the sequence may not be captured, which may lead to sub-optimal solutions. A possible remedy for this is to define the randomization tasks as sets of different and possibly intersecting groups of randomization parameters [41]. Furthermore, this paper does not address how to define the ranges for the randomization parameters. Excessive randomization is possible even when a single parameter is randomized at a time, e.g., 10% noise randomization translates to ± 200 mm uncertainty of the box's position and makes the grasping task challenging to solve, leading to poor grasp success on the real system. To address this issue, CDR could be combined with *automated domain randomization* [14] or *active domain randomization* [15] to find suitable ranges for each randomization parameter.

In summary, CDR provides a flexible framework for zero-shot sim2real transfer that does not require all randomization parameters to be defined or implemented ahead of time. It enables randomization parameter decoupling and continual model adaptation on new randomizations if required. In future work, we want to address the limitations mentioned above, as well as to implement CDR with other RL and regularization-based CL approaches and apply it in the context of vision-based domain randomization.

REFERENCES

- [1] W. Zhao, J. P. Queralta, and T. Westerlund, “Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: A Survey,” in *IEEE Symposium Series on Computational Intelligence (SSCI)*, Canberra, ACT, Australia, Dec. 2020, pp. 737–744.
- [2] G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester, “Challenges of Real-World Reinforcement Learning: Definitions, Benchmarks and Analysis,” *Machine Learning*, vol. 110, no. 9, pp. 2419–2468, 2021.
- [3] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, “Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 23–30.
- [4] J. Josifovski, M. Malmir, N. Klarmann, B. L. Zagar, N. Navarro-Guerrero, and A. Knoll, “Analysis of Randomization Effects on Sim2Real Transfer in Reinforcement Learning for Robotic Manipulation Tasks,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Kyoto, Japan: IEEE, Oct. 2022, pp. 10 193–10 200.
- [5] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, “Continual Lifelong Learning with Neural Networks: A Review,” *Neural Networks*, vol. 113, pp. 54–71, May 2019.
- [6] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal Policy Optimization Algorithms,” Tech. Rep. arXiv: 1707.06347, July 2017.
- [7] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, “Overcoming Catastrophic Forgetting in Neural Networks,” *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [8] J. Schwarz, W. Czarnecki, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R. Hadsell, “Progress & Compress: A Scalable Framework for Continual Learning,” in *International Conference on Machine Learning (ICML)*, vol. 80. PMLR, 2018, pp. 4528–4537.
- [9] S. James and E. Johns, “3D Simulation for Robot Arm Control with Deep Q-Learning,” in *NIPS Workshop: Deep Learning for Action and Interaction*, Barcelona, Spain, Dec. 2016.
- [10] M. Kaspar, J. D. Muñoz Osorio, and J. Bock, “Sim2Real Transfer for Reinforcement Learning Without Dynamics Randomization,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Las Vegas, NV, USA: IEEE, Oct. 2020, pp. 4383–4388.
- [11] F. Muratore, F. Ramos, G. Turk, W. Yu, M. Gienger, and J. Peters, “Robot Learning from Randomized Simulations: A Review,” *Frontiers in Robotics and AI*, vol. 9, no. 799893, Apr. 2022.
- [12] J. Josifovski, M. Kerzel, C. Pregizer, L. Posniak, and S. Wermter, “Object Detection and Pose Estimation Based on Convolutional Neural Networks Trained with Synthetic Data,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Madrid, Spain, Oct. 2018, pp. 6269–6276.
- [13] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-Real Transfer of Robotic Control with Dynamics Randomization,” in *IEEE International Conference on Robotics and Automation (ICRA)*, Brisbane, QLD, Australia, 2018, pp. 3803–3810, iSSN: 2577-087X.
- [14] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, “Solving Rubik’s Cube with a Robot Hand,” Tech. Rep. arXiv:1910.07113, Oct. 2019.
- [15] B. Mehta, M. Diaz, F. Golemo, C. J. Pal, and L. Paull, “Active Domain Randomization,” in *Conference on Robot Learning (CoRL)*, vol. 100. Osaka, Japan: PMLR, May 2020, pp. 1162–1176.
- [16] F. Ramos, R. Possas, and D. Fox, “BayesSim: Adaptive Domain Randomization Via Probabilistic Inference for Robotics Simulators,” in *Robotics: Science and Systems (R:SS)*, vol. 15, June 2019.
- [17] F. Muratore, T. Gruner, F. Wiese, B. Belousov, M. Gienger, and J. Peters, “Neural Posterior Domain Randomization,” in *Conference on Robot Learning (CoRL)*, vol. 164. PMLR, Nov. 2021, pp. 1532–1542.
- [18] P. Huang, X. Zhang, Z. Cao, S. Liu, M. Xu, W. Ding, J. Francis, B. Chen, and D. Zhao, “What Went Wrong? Closing the Sim-to-Real Gap Via Differentiable Causal Discovery,” in *Conference on Robot Learning (CoRL)*, Atlanta, GA, USA, Aug. 2023.
- [19] G. Tiboni, K. Arndt, and V. Kyrki, “DROPO: Sim-to-Real Transfer with Offline Domain Randomization,” *Robotics and Autonomous Systems*, vol. 166, p. 104432, Aug. 2023.
- [20] P. Schöpf, S. Auddy, J. Hollenstein, and A. Rodríguez-Sánchez, “Hypernetwork-PPO for Continual Reinforcement Learning,” in *Deep Reinforcement Learning Workshop NeurIPS 2022*, Dec. 2022.
- [21] S. Auddy, J. Hollenstein, M. Saveriano, A. Rodríguez-Sánchez, and J. Piater, “Continual Learning from Demonstration of Robotics Skills,” *Robotics and Autonomous Systems*, vol. 165, p. 104427, 2023.
- [22] —, “Scalable and Efficient Continual Learning from Demonstration Via a Hypernetwork-Generated Stable Dynamics Model,” Tech. Rep. arXiv:2311.03600, Jan. 2024, eprint: 2311.03600.
- [23] K. Khetarpal, M. Riemer, I. Rish, and D. Precup, “Towards Continual Reinforcement Learning: A Review and Perspectives,” *Journal of Artificial Intelligence Research*, vol. 75, pp. 1401–1476, Dec. 2022.
- [24] Y. Luo, Y. Wang, K. Dong, Q. Zhang, E. Cheng, Z. Sun, and B. Song, “Relay Hindsight Experience Replay: Self-Guided Continual Reinforcement Learning for Sequential Object Manipulation Tasks with Sparse Rewards,” *Neurocomputing*, vol. 557, p. 126620, Nov. 2023.
- [25] J. Josifovski, M. Malmir, N. Klarmann, and A. Knoll, “Continual Learning on Incremental Simulations for Real-World Robotic Manipulation Tasks,” in *2nd R:SS Workshop on Closing the Reality Gap in Sim2Real Transfer for Robotics*, Corvallis, OR, USA, July 2020, p. 3.
- [26] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, “Progressive Neural Networks,” Tech. Rep. arXiv:1606.04671, Oct. 2016.
- [27] A. A. Rusu, M. Veer, T. Rothrl, N. Heess, R. Pascanu, and R. Hadsell, “Sim-to-Real Robot Learning from Pixels with Progressive Nets,” in *Annual Conference on Robot Learning (CoRL)*, vol. 78. Mountain View, CA, USA: PMLR, 2017, pp. 262–270, iSSN: 2640-3498.
- [28] A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirkpatrick, R. Pascanu, V. Mnih, K. Kavukcuoglu, and R. Hadsell, “Policy Distillation,” in *International Conference on Learning Representations (ICLR)*. San Juan, Puerto Rico: arXiv, May 2016, eprint: 1511.06295.
- [29] R. Traoré, H. Caselles-Dupré, T. Lesort, T. Sun, N. Díaz-Rodríguez, and D. Filliat, “Continual Reinforcement Learning Deployed in Real-Life Using Policy Distillation and Sim2real Transfer,” in *ICML Workshop on Multi-Task and Lifelong Learning*. arXiv, June 2019.
- [30] S. Kessler, J. Parker-Holder, P. Ball, S. Zohren, and S. J. Roberts, “UNCLEAR: A Straightforward Method for Continual Reinforcement Learning,” in *ICML Workshop on Continual Learning*, vol. 108. Vienna, Austria: PMLR, 2020.
- [31] A. Gaurav, “Safety-Oriented Stability Biases for Continual Learning,” Master’s thesis, University of Waterloo, 2020.
- [32] “Robotiq Gripper.” [Online]. Available: <http://robotiq.com/products/industrial-robot-hand/>
- [33] P. Varin, L. Grossman, and S. Kuindersma, “A Comparison of Action Spaces for Learning Manipulation Tasks,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Macau, China, 2019, pp. 6015–6021.
- [34] Y. R. Stürz, L. M. Afholter, and R. S. Smith, “Parameter Identification of the KUKA LBR iiwa Robot Including Constraints on Physical Feasibility,” *IFAC-PapersOnLine*, vol. 50, no. 1, pp. 6863–6868, 2017.
- [35] S. Edwards and C. Lewis, “ROS-Industrial — Applying the Robot Operating System (ROS) to Industrial Applications,” in *ECHORD Workshop at the IEEE International Conference on Robotics and Automation (ICRA)*, St. Paul, MN, USA, 2012.
- [36] C. Hennersperger, B. Fuerst, S. Virga, O. Zettinig, B. Frisch, T. Neff, and N. Navab, “Towards MRI-Based Autonomous Robotic Us Acquisitions: A First Feasibility Study,” *IEEE Transactions on Medical Imaging*, vol. 36, no. 2, pp. 538–548, Feb. 2017.
- [37] M. Quigley, B. Gerkey, K. Conley, J. Faust, T. Foote, J. Leibs, E. Berger, R. Wheeler, and A. Ng, “ROS: An Open-Source Robot Operating System,” in *ICRA Workshop on Open Source Software*, vol. 3, 2009, p. 6.
- [38] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu, “Stable Baselines,” 2018.
- [39] A. Raffin, J. Kober, and F. Stulp, “Smooth Exploration for Robotic Reinforcement Learning,” in *Conference on Robot Learning (CoRL)*, vol. 164. London, UK: PMLR, Jan. 2022, pp. 1634–1644.
- [40] M. Malmir, J. Josifovski, N. Klarmann, and A. Knoll, “DiAREL: Reinforcement Learning with Disturbance Awareness for Robust Sim2Real Policy Transfer in Robot Control,” Tech. Rep. arXiv:2306.09010, 2023.
- [41] Y. Kadokawa, L. Zhu, Y. Tsurumine, and T. Matsubara, “Cyclic policy distillation: Sample-efficient sim-to-real reinforcement learning with domain randomization,” *Robotics and Autonomous Systems*, vol. 165, p. 104425, 2023.