

---

# Universal benchmark for actuation dynamics adaptation in reinforcement learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Enabling reinforcement learning (RL) agents to adapt to changing environment  
2 dynamics is crucial for robustness. Consider the case where a robot’s motors and  
3 gears change their behavior due to wear and tear over time, or where an old used-  
4 up part gets replaced. Current literature primarily emphasizes resilience against  
5 observation noise, distractions in the environment or shifts in physical properties  
6 of the world. However, the problem of continual shifts or sudden changes in  
7 actuation dynamics is relatively unexplored. To facilitate systematic research  
8 in that regard, we contribute a *Universal Benchmark for Actuation Dynamics  
9 Adaptation (UBADA)*<sup>1</sup>. We present a universal set of wrappers compliant with  
10 the Gymnasium API standard, providing a multitude of challenges with continual  
11 (serial) and multi-task (parallel) learning scenarios of changing action dynamics.  
12 We showcase the problem on visual and low-dimensional proprioceptive inputs,  
13 with dense or sparse rewards, utilizing the state-of-the-art learning algorithms  
14 Soft-Actor-Critic (SAC) and Data-regularized Q (DrQ).

## 15 1 Introduction

16 In current research to systematically investigate transfer capabilities of a reinforcement learning (RL)  
17 agent, one typically asks questions like the following: "The agent can *open a door*, but can it *open a  
18 window*?" [62, 61, 50, 47, 58]. Benchmarks like Meta-World [62] or CompoSuite [37] contain 10, 50  
19 or up to 256 different tasks like those used in the example above. Another perspective on transfer  
20 learning involves exploring generalization and robustness to perturbations in the observation [18].  
21 For example, one could ask: "The agent can *open a door*, but can it also *open the door if the floor is  
22 distracted with colorful dots*?".

23 With our contribution, we emphasize an often neglected angle on robustness and the ability to adapt  
24 to changes in the agent’s embodiment. For example, we would ask: "The agent can *open a door*, but  
25 could it also *open a door if it were weaker*?". To address such questions, we change the dynamics of  
26 the system by manipulating the agent’s action effect.

27 To illustrate, assuming you are reading this work on a computer. Navigate to your system settings and  
28 invert the scrolling direction of your trackpad or mouse. With the switched scrolling direction, return  
29 to the document and resume reading. You will notice that adjusting to the new system behavior might  
30 require a few trials, but overall you can adapt relatively quickly, and you do not need to re-learn to  
31 operate the trackpad or mouse from scratch.

32 A possible motor psychological explanation for how humans adapt to new dynamics involves a  
33 shift of attention to the relevant sensory modalities. Humans have efficient innate abilities to handle

---

<sup>1</sup><https://github.com/xxx/xxx>

34 conflicting sensory stimuli. When using the inverted trackpad or computer mouse, humans suppress  
35 the attention on haptic and proprioceptive senses and focus on the unfamiliar visual stimulus on  
36 the screen. The unfamiliar visual stimuli are sensory conflicts that emerge when novel experiences  
37 contradict expectations that have been formed through lifelong learning [32]. See Figure 1 for an  
38 illustration of the inverted computer mouse and Section A in the Appendix for more examples.

39 The example with the inverted mouse is rather  
40 constructed, but it is a good proxy for many re-  
41 alistic cases in robotics. For example, wear and  
42 tear of motors can lead to fluctuations in their  
43 effectiveness over time. Replacing a motor or a  
44 joint yields different dynamics that have to be  
45 relearned. Even in the case of a malfunction-  
46 ing part, the robot should be able to adapt and  
47 compensate for the missing functionality.

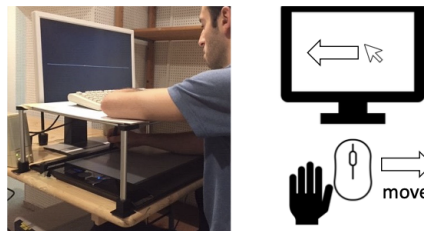


Figure 1: Psychological experiments requiring subjects to learn novel unexpected dynamics [13, 31, 32], e.g., moving the hand to the right causes the mouse cursor to move left. Illustration adapted from Liesner et al. [32].

48 In order to address such problems or similar  
49 ones in reinforcement learning (RL), significant  
50 effort is being devoted to improving robustness,  
51 (zero-shot) adaptation and overall generalization  
52 [59, 18, 49, 24]. However, current RL bench-  
53 marks and methods primarily focus on address-  
54 ing noise, augmentation, perturbation, or distrac-  
55 tions at the observation level, often neglecting  
56 changes in actuation dynamics.

57 To bring greater attention to this aspect, we propose our main contribution: a *Universal Benchmark*  
58 *for Actuation Dynamics Adaptation (UBADA)*, designed to analyze the adaptation capabilities of RL  
59 agents when confronted with dynamic changes in the agent’s actuation. UBADA can be considered  
60 a set of universal wrappers for environments following the Gymnasium API standard [53]. It  
61 is *universal*, in the sense that it is not limited to specific hand-crafted environments, and it is  
62 independent of the type of sensory input. Therefore, it can convert a variety of publicly available  
63 standard environments into a challenging task for the continual learning, curriculum learning (serial  
64 transfer) or multi-task learning (parallel transfer) setup. In addition, it can be used to extend existing  
65 benchmarks by increasing the number of available tasks. The focus on adaptation and robustness  
66 and the universal nature of our benchmark facilitates research into AutoRL techniques aimed at  
67 optimizing agents for generality and applicability [4, 38]. The second contribution of this work  
68 involves evaluating the transfer capabilities of two state-of-the-art RL algorithms, namely SAC  
69 [16] and DrQ [60]. Therefore, we utilize environments from Gymnasium MuJoCo [53, 52] (goal-  
70 conditioned and non goal-conditioned) and DM Control [54].

## 71 2 Related work

72 Our contribution is a benchmark for changing action dynamics. Therefore, it is related to transfer  
73 learning, multi-task and meta learning, continual and curriculum learning. Another relevant field of  
74 research is RL in modified environments.

### 75 2.1 Transfer, multi-task and continual learning benchmarks

76 Henderson et al. [19] regard the Arcade Learning Environment (ALE) [3] as the primary benchmark  
77 for evaluating multi-task learning in domains with discrete actions. However, for continuous actions  
78 they observe a lack of a standardized evaluation environments for multi-task learning at that time.  
79 To fill this gap, they introduce a benchmark [19], aiming to facilitate a systematic comparison of  
80 multi-task, transfer, and lifelong learning in continuous domains.

81 The creators of Meta-World [62] emphasize the possible occurrence of negative transfer between  
82 tasks in the ALE. Hence, they motivate their work by the idea that positive transfer between different  
83 tasks should be possible. They propose a set of related yet diverse robotics tasks that share the same  
84 robot, action space, and workspace. It became a very popular benchmark used to study multi-task  
85 and meta learning [61, 63, 47, 50].

86 In robotics, more benchmarks comprising a multitude of environments exist [14, 22, 54, 58, 37]. For  
87 example, CompoSuite [37] introduces a series of tasks where a robotic arm manipulates individual  
88 objects to achieve objectives while navigating obstacles. This framework facilitates the evaluation of  
89 RL approaches in learning the compositional structure of tasks and their ability for compositional  
90 generalization to unseen tasks. Avalanche RL [35, 5] provides support for a continuous stream  
91 of diverse environments. Additionally, CORA [40] encompasses metrics and baselines designed  
92 to assess various aspects of continual RL. These aspects include catastrophic forgetting, plasticity,  
93 generalization, and sample-efficient learning. The benchmarks are based on diverse environments  
94 such as ALE [3, 25, 45, 43], ProcGen [8, 21], NetHack [28, 44], ALFRED [46] (built upon AI2-  
95 THOR [26]). While many efforts in continual learning tend to emphasize catastrophic forgetting,  
96 Continual World [58] advocates for the significance of forward transfer. Another notable contribution  
97 is BabyAI [7], which involves the creation of multiple grid-world environments of increasing difficulty.  
98 Here, the primary focus is on investigating grounded language learning.

## 99 2.2 Environment modification benchmarks

100 The perceived dynamics of an environment can be changed by directly modifying the agent’s  
101 observation. For tasks derived from pixel inputs, this could involve the incorporation of color and  
102 background distractions as valuable benchmarks [18, 49]. Another approach is a straightforward data  
103 augmentation technique, as proposed by Yarats et al. [60]. The authors utilize perturbations on the  
104 input observations commonly employed in computer vision tasks to regularize the value function.  
105 With DrQ they show, that this augmentation approach proves to be particularly effective in enhancing  
106 the performance of the SAC algorithm [16]. Not necessarily focused on RL from pixels, in their  
107 study, Khraishi and Okhrati [23] also delve into the application of data augmentation techniques in  
108 RL to enhance model performance and promote diversity in training data.

109 Another way to modify the perceived dynamics of an environment is by using action noise. Action  
110 noise plays an important role for exploration behavior [20, 12], as in Deep Deterministic Policy  
111 Gradient (DDPG) [33] and Twin Delayed DDPG (TD3) [15]. It is commonly generated by utilizing  
112 the Ornstein–Uhlenbeck process [55] or drawing from uncorrelated Gaussian distributions. In the  
113 case of SAC [16], where the policy is stochastic, Gaussian noise is implicitly introduced through the  
114 sampling procedure. In our work, the modification of actuation dynamics serves a different purpose  
115 than exploration: We focus on analyzing the adaptability of agents to dynamic changes.

116 Langlois and Everitt [29] have the human-in-a-loop aspect in mind and analyze how an agent behaves  
117 if its intended actions are overridden or manipulated by the environment or another agent. That is  
118 also related to the work on safety in RL by Leike et al. [30].

119 van Seijen et al. [56] want to measure how quickly an agent that is trained on task A changes its  
120 policy after it is placed in task B. They specifically experiment with maintaining identical transition  
121 dynamics between tasks A and B, while introducing a local difference in the reward function. For  
122 that purpose they introduce the local change adaptation (LoCA) regret metric.

123 Similar to our method, the CARL benchmark makes environments configurable and facilitates training  
124 agents to generalize across different instances (contexts) of the same environment. This supports  
125 research into AutoRL techniques aimed at optimizing agents for broad generality [4, 38].

126 Furthermore, there is a trend in research to modify or perturb system parameters, such as body part  
127 size or gravity, in continuous control tasks [19, 11].

128 The work that is most closely related to ours is the Real-World Reinforcement Learning (RWRL)  
129 Challenge Framework [11]. It provides a diverse suite of tasks to change the morphology of the  
130 agent or the physics of the world, like the friction of the ground. Related to the action effect it allows  
131 action offsets, action noise, repetitive actions, and action drops. However, it exclusively includes  
132 proprioceptive observations and does not allow for visual observations, and it is constrained by a  
133 limited, predefined set of environments. While the concepts from RWRL could in principle be applied  
134 to novel environments, doing so within the RWRL system may pose considerable challenges. In  
135 our approach, we address these drawbacks to provide a more general solution which is implemented  
136 as a set of universal wrappers that modify actions. Furthermore, our approach facilitates extensive  
137 configurability. For instance, action effects can be altered selectively for individual action dimensions,  
138 e.g., specific joints in a robot arm, rather than modifying the action as a whole.

### 139 3 Background

140 Reinforcement learning builds on Markov Decision Processes (MDP), defined as tuples  
141  $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$ . The behavior of the world is determined by the transition probability  $P(s'|s, a)$ ,  
142 whereas the policy  $\pi(a|s)$  generates an action  $a \in \mathcal{A}$  based on inputs  $s \in \mathcal{S}$ , to maximize the cumula-  
143 tive return, the sum of discounted rewards  $\sum_{t=0} \gamma^t r_t$ , with  $r_t = R(s_t, a_t, s_{t+1})$ . The particularity of  
144 this work is that we introduce a modification of the action  $p(\bar{a}|a)$  which changes over time. Hence, in  
145 the transition dynamics  $P(s'|s, a)$  the action is manipulated by our wrappers to yield  $P(s'|s, \bar{a})$ . The  
146 outcome could be similar or totally diverge from the old one.

147 An extended background section on transfer learning in RL, multi-task RL and continual RL is  
148 provided in the Appendix B.

### 149 4 Manipulate the action effect

150 UBADA is designed around the Gymnasium API standard [53]. It consists of a variety of wrappers  
151 that manipulate the action effect of the agent within its environment. They are universally applicable  
152 insofar as they are independent of the selected base environment and independent of whether pixels  
153 or state features are processed. Regardless of the action effect wrapper selected, the action space  
154 remains the same as in the original base environment. The wrappers provide versatility by allowing  
155 toggling on and off, seamless combination with each other, and the capability to alter the action effect  
156 of one or all action dimensions. Additionally, they are configurable. This design grants considerable  
157 freedom in tailoring the experimental setup according to specific preferences.

158 Note that the possible modifications for discrete actions are limited due to their inherent limited  
159 flexibility, the focus is on continuous actions.

160 The action effect wrappers are listed in the following. An overview including a brief description and  
161 motivation can be found in Table 1 in the Appendix. In general, all wrappers are configurable by the  
162 choice of the action dimension which is modified and if applicable by another specific parameter  
163 which is mentioned in the following.

164 **InvertAction** Consider the concept of the backward bicycle, where attempting a left turn of the  
165 handlebar leads to movement in the opposite, or right, direction (cf. Appendix A). This wrapper  
166 inverts the continuous action, simply by switching the sign of one or all action dimensions.

167 **ScaleAction** This wrapper scales the continuous action by multiplying either one or all action  
168 dimensions with a scalar value. A scenario where this might be relevant is to simulate decaying motor  
169 efficiencies over a robot’s lifetime. The scaling factor is configurable.

170 **OffsetAction** Similarly, this wrapper adds a constant scalar value to one or all dimensions of the  
171 continuous action. This might simulate an inadequate calibration of a system, reflecting a systematic  
172 error. The offset is configurable.

173 **NoiseAction** This wrapper adds a random value to either one or all dimensions of the continuous  
174 action. The random value is drawn from a Gaussian distribution with a mean  $\mu = 0$  and a user-defined  
175 standard deviation  $\sigma$ . While being applied at the environment level, it can be perceived as analogous  
176 to Gaussian action noise commonly utilized in algorithms like DDPG or TD3. In a robotics context,  
177 this could simulate the presence of random errors. The choice of  $\sigma$  is configurable.

178 **SineNoiseAction** Expanding upon the **NoiseAction** wrapper, this implementation adds a sine  
179 offset to the Gaussian noise. A single parameter within the wrapper specifies both the standard  
180 deviation  $\sigma$  and the amplitude of the sinusoidal wave. This allows for analyzing the impact of both an  
181 unpredictable component (Gaussian noise) and a predictable element (sine offset) on the action. The  
182 choice of  $\sigma$  is configurable.

183 **ZeroAction** It randomly sets one or all dimensions of the continuous action to zero for a defined  
184 number of steps. Again, in a robotics context, this wrapper might correspond to loose contacts or  
185 sporadic complete engine failures. The probability of a dimension (or the whole action) being zeroed  
186 and the duration of this zeroing effect are both configurable.

187 **RepeatAction** Similar to **ZeroAction**, but instead of setting to zero this wrapper repeats one  
188 or all dimensions of the continuous or discrete action for a defined number of steps. Proposed by

189 Machado et al. [36] this effect is also known as sticky actions, though it is usually not possible to  
190 repeat only one dimension of the action. The probability of a dimension (or the whole action) being  
191 repeated and the duration of this repetition are both configurable.

192 `SwapAction` This wrapper swaps either one dimension of the continuous or discrete action with  
193 another randomly picked one or shuffles randomly all dimensions. The changed order of dimensions  
194 is kept for the whole period of time this wrapper is activated. Naively, one can imagine interchanged  
195 cables wrongly connecting the motors in a robot system.

196 An additional `DynamicsHintObservation` wrapper can be applied optionally, it augments the  
197 observation space with a one-hot encoding which functions as a task identifier and informs the agent  
198 about the applied and active action effect wrapper.

## 199 5 Action effect benchmark

200 In the context of RL, the term *task* can sometimes be subject to varying interpretations. In our  
201 definition, each dynamic modification is considered a distinct task, or context, even though the  
202 underlying environment remains constant.

203 In this section we provide experimental results using two environments: `HalfCheetah-v4` based on  
204 work by Wawrzyński [57] and available through Towers et al. [53] and `walker-walk-v0` from DM  
205 Control [54], utilized via Towers et al. [53] and Tai et al. [51] with proprioceptive inputs. Extended  
206 results for more kinds of environments (goal-conditioned, visual observation inputs) showcasing the  
207 versatility of the wrapper approach are provided in Appendix G (cf. Table 3 for an overview). We use  
208 the learning algorithms SAC [16] for proprioceptive and DrQ [60] for visual input observations. For  
209 goal-conditioned environments we combine SAC with Hindsight Experience Replay (HER) [1]. We  
210 use implementations provided by Kostrikov [27] (JAXRL2) and Raffin et al. [42] (Stable-Baselines  
211 3). Hyperparameters are kept default as provided in their implementations (cf. Section D in the  
212 Appendix). In general, results are obtained by averaging across five runs with different random  
213 initialization.

214 On the one hand, our intent is to demonstrate transfer capabilities of general purpose agents while  
215 learning under changing conditions. Although very interesting, we do not investigate how state-  
216 of-the-art approaches specifically designed for multi-task RL [61, 47, 50] or continual RL [25, 43]  
217 behave under action dynamic changes. We also do not analyze the effect of different replay buffer  
218 sizes which possibly influences learning and transfer capabilities in the continual learning case.

219 On the other hand, the purpose of the following experiments and those in Section G in the Appendix  
220 is to illustrate the universal applicability of action modifications across various environments. In some  
221 cases the wrappers yield fundamentally different modifications. Not all modifications are supposed to  
222 be transferable to real-world scenarios. Still, they can be useful to reveal the limits of the agent, in  
223 particular with regard to dimensions of robustness and transfer capabilities. An extended discussion  
224 on the classification of the different modifications is provided in Section E in the Appendix.

### 225 5.1 Sequential training setup (continual learning)

#### 226 5.1.1 Experimental setting

227 The use of wrappers on top of the actual environments functions like a modular system. Changes  
228 can be combined with each other as desired and triggered at any time step. For the sequential  
229 experiments, we first use the initial environment and then switch the dynamic of the environment by  
230 toggling the respective action effect wrapper after a defined number of steps. For the environments  
231 `walker-walk-v0` and `HalfCheetah-v4` we train each task  $\mathcal{T}_i$  for  $T_{\mathcal{T}} = 500000$  and  $T_{\mathcal{T}} = 1000000$   
232 steps, respectively. The choice of  $T_{\mathcal{T}}$  is made to ensure that, at the very least, the unmodified  
233 environment can be trained to a reasonable degree. We only consider sequences where  $T_{\mathcal{T}}$  is  
234 constant across all tasks. For example, if we consider the initial dynamic and nine modifications  
235 of the `walker-walk-v0` environment, the sequence has  $N = 10$  tasks and in total we train for  
236  $T = N \cdot T_{\mathcal{T}} = 5000000$  steps. The  $i$ -th task is trained during the interval  $t \in [(i - 1) \cdot T_{\mathcal{T}}, i \cdot T_{\mathcal{T}}]$ .

237 Unlike Powers et al. [40] we do not cycle through the sequence of tasks multiple times. (Powers et al.  
238 [40] also only cycle multiple times in the general case, however, to compute forward transfer and  
239 forgetting they use one cycle.)

240 While it is possible to have multiple wrappers activated simultaneously, for a clearer understanding  
 241 of the impact of the individual wrappers, we focus on a scenario where, at each time step, only one  
 242 wrapper is active. Moreover, within our experiments we do not combine multiple types of wrappers,  
 243 e.g., `ScaleAction` and `OffsetAction`. Possibilities are extensive.

### 244 5.1.2 Metrics

245 The metrics considered and their implementation were mainly based on the work of Lopez-Paz and  
 246 Ranzato [34], Díaz-Rodríguez et al. [10], Wołczyk et al. [58] and Powers et al. [40].

247 Let  $R_t(\pi, \mathcal{T}_i)$  be the episodic return, the undiscounted sum of rewards received over an episode,  
 248 under the policy  $\pi$  on task  $\mathcal{T}_i$  at time step  $t$ . To enhance the understanding of results and to account  
 249 for different attainable returns in different environments we define a performance measure  $\rho_i(t)$  by  
 250 normalizing the return by its maximum value over all time steps and over all tasks. Note, within a  
 251 sequence each task is derived from the same environment, hence we use the maximum over all tasks  
 252 which differs from Powers et al. [40]. When comparing performances of two tasks  $\mathcal{T}_i$  and  $\mathcal{T}_j$  within  
 253 one task sequence we generally assume  $i < j$ .

$$\rho_i(t) = \frac{R_t(\pi, \mathcal{T}_i)}{\max_{1 \leq t \leq T, 1 \leq k \leq N} R_t(\pi, \mathcal{T}_k)} \quad (1)$$

254 **Average performance** The performance at any time step averaged across all tasks. We consider its  
 255 final value  $P(t = T)$  for evaluation.

$$P(t) = \frac{1}{N} \sum_{i=1}^N \rho_i(t) \quad (2)$$

256 **Forward transfer** Measures the influence that learning a task  $\mathcal{T}_i$  has on the performance of a  
 257 future task  $\mathcal{T}_j$ . It can occur when the model is able to perform zero-shot learning [10], e.g., Powers  
 258 et al. [40] refer to the corresponding metric as zero-shot forward transfer.

$$FT = \frac{2}{N(N-1)} \sum_{j=2}^N \sum_{i=1}^{j-1} FT_j \quad \text{where} \quad FT_j = \rho_j(t = i \cdot T_{\mathcal{T}}) - \rho_j(t = (i-1) \cdot T_{\mathcal{T}}) \quad (3)$$

259 **Backward transfer** In contrast to the findings reported by Wołczyk et al. [58], where they observe  
 260 no occurrence of backward transfer in their scenario, we find the prospect intriguing for our specific  
 261 purposes. Differences between tasks in our context may be more nuanced, implying that learning  
 262 a modification of an environment could potentially enhance performance on the initial task in a  
 263 retrospective manner.

264 To be precise, we actually use positive backward transfer following the argumentation of Díaz-  
 265 Rodríguez et al. [10] and Wołczyk et al. [58].

$$BT = \frac{2}{N(N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N BT_i \quad \text{where} \quad BT_i = \max\{0, \rho_i(t = j \cdot T_{\mathcal{T}}) - \rho_i(t = (j-1) \cdot T_{\mathcal{T}})\} \quad (4)$$

266 **Forgetting** It represents the decrease of performance from a learned task during later tasks.

$$F = \frac{2}{N(N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N F_i \quad \text{where} \quad F_i = \max\{0, \rho_i(t = (j-1) \cdot T_{\mathcal{T}}) - \rho_i(t = j \cdot T_{\mathcal{T}})\} \quad (5)$$

267 In alignment with the argumentation presented in Powers et al. [40] and akin to the approach taken  
 268 by Lopez-Paz and Ranzato [34], unlike Chaudhry et al. [6] we use the performance right before

269 training on a new task rather than the maximum value observed so far. That approach corresponds to  
 270 a more task isolated view, hence Powers et al. [40] denotes their metric isolated forgetting. Wołczyk  
 271 et al. [58] consider the decrease of performance after ending the training on the whole task sequence.  
 272 They compute  $F_i = \rho_i(t = i \cdot T_{\mathcal{T}}) - \rho_i(t = T)$ , which corresponds to a more cumulative effect of  
 273 forgetting.

274 Also note, to emphasize two different concepts, we explicitly distinguish positive backward transfer  
 275 and forgetting, although Equation 4 and 5 are quite related [10].

### 276 5.1.3 One dynamic switch

277 In this scenario, a SAC agent trains for  $T_{\mathcal{T}} = T/2$  steps on a task 1,  $\mathcal{T}_1$ , and continues training on  
 278 task 2,  $\mathcal{T}_2$ , for the same number of steps.  $\mathcal{T}_1$  is the unmodified environment.  $\mathcal{T}_2$  comes from the same  
 279 environment, but with modified action dynamics. The modification is defined by the type of action  
 280 wrapper, the action dimension to be modified, and if applicable a wrapper specific configuration  
 281 value. In Figure 2 we use the `InvertAction` and the `ScaleAction` wrapper. We modify only action  
 282 dimension 0. For `ScaleAction` we consider three scaling factors  $\{0.2, 0.5, 0.8\}$ .

283 With the one dynamic switch setup and the selection of action effect modifications we showcase the  
 284 different metrics defined in Section 5.1.2. While the `InvertAction` switch is a perfect example  
 285 for (catastrophic) forgetting ( $F = 0.92$ ), we cannot observe any forgetting of the first task after the  
 286 `ScaleAction` switch (cf. Figure 2a). It applies analogously for the forward transfer. Just considering  
 287 individual task returns, the agent is not able to transfer knowledge from the baseline environment  
 288 to the same environment with the `InvertAction` modification. `ScaleAction` allows for forward  
 289 transfer, correlating with the manifestation of the switch. The lower the scaling value, the higher the  
 290 contrast between tasks. For scaling values 0.8, 0.5 and 0.2 we have  $FT = 0.96$ ,  $FT = 0.89$  and  
 291  $FT = 0.49$ , respectively. To analyze positive backward transfer we consider the `HalfCheetah-v4`  
 292 environment (cf. Figure 2b). For a `ScaleAction` switch and lower scaling values we can observe  
 293 positive backward transfer, while forgetting emerges otherwise. For scaling values 0.8, 0.5 and 0.2  
 294 we have  $BT = 0.12$ ,  $BT = 0.02$ , and  $BT = 0.0$ , respectively.

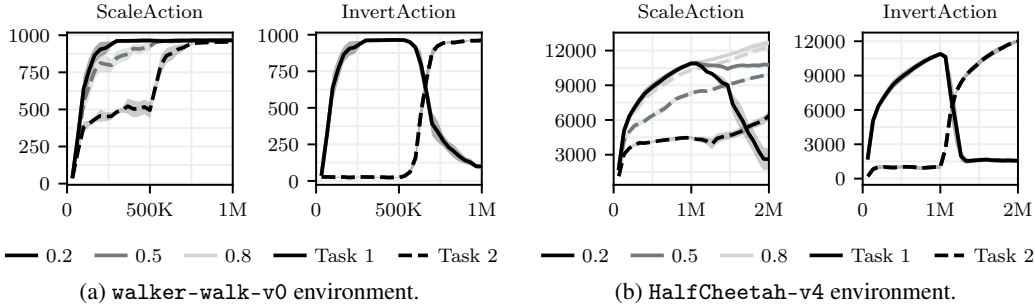


Figure 2: Per task evaluation returns for the sequential training setup with one dynamic switch. To establish the switch we use the `InvertAction` and the `ScaleAction` wrapper. We modify only action dimension 0. For `ScaleAction` we consider three scaling factors  $\{0.2, 0.5, 0.8\}$ . The switch happens at  $T/2$ .

### 295 5.1.4 Continual adaptation to new dynamics

296 In this scenario a SAC agent trains for  $T_{\mathcal{T}} = T/10$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ .  
 297 Then nine times for  $T_{\mathcal{T}} = T/10$  steps on modified action dynamics, varied by the wrapper value,  
 298 tasks  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{10}\}$ . We construct the sequence of tasks with increasing difficulty. One could  
 299 think of a curriculum learning schedule. For `NoiseAction` tasks are determined by the standard  
 300 deviation for the Gaussian noise in  $\{0.1, 0.2, \dots, 0.9\}$ , for `OffsetAction` the action offset is in  
 301  $\{-0.1, -0.2, \dots, -0.9\}$  and for `ScaleAction` the scaling factor is in  $\{0.9, 0.8, \dots, 0.1\}$ .

302 Consider solely the training return in Figure 3a, the agent can adapt to the action effect changes  
 303 smoothly without larger collapses. In that regard there is no difference between the three action  
 304 effect wrapper. Now consider the evaluation returns for each individual task at every time step in  
 305 Figure 3b. Differences in forward transfer and forgetting stand out. While the forward transfer in

306 the NoiseAction and OffsetAction experiments mostly happens during the training of task  $\mathcal{T}_1$ ,  
 307 in the ScaleAction experiment it occurs on a wider time span, i.e., during the training of  $\mathcal{T}_1$  to  $\mathcal{T}_5$ .  
 308 Nevertheless, overall according to Equation 3 all three experiments have the same total  $FT = 0.19$ .  
 309 As the metric is averaged across tasks and training periods, it cannot provide an isolated view on  
 310 forward transfer. Observing Figure 3b or the distribution of  $FT_j$  (cf. Equation 3) is more appropriate.  
 311 Towards the end of the training sequence the agent shows to some degree forgetting of the earlier  
 312 tasks in the context of the OffsetAction experiment. The longer it is been since the training of  
 313 a task, the more is forgotten. Interestingly, this notion is not observed in the NoiseAction and  
 314 ScaleAction experiments.

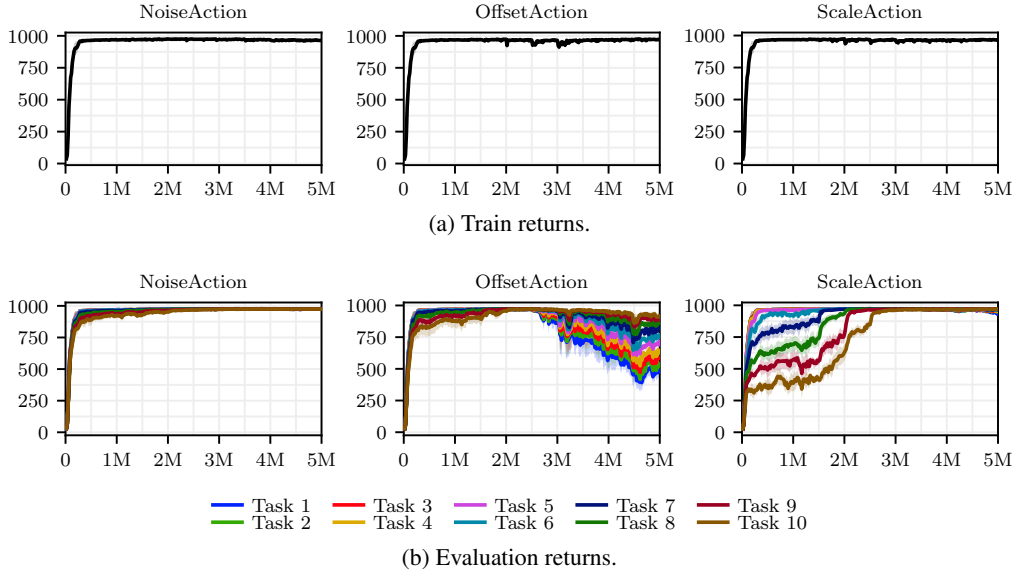


Figure 3: Returns for the sequential training setup with continual adaptation. A SAC agent is trained on the walker-walk-v0 environment. After first training on the unmodified action dynamic, action effect changes occur every  $T/10$  steps. For NoiseAction tasks are determined by an action wrapper value in  $\{0.1, 0.2, \dots, 0.9\}$ , for OffsetAction in  $\{-0.1, -0.2, \dots, -0.9\}$ , for ScaleAction in  $\{0.9, 0.8, \dots, 0.1\}$ . The modification is done on action dimension 0.

## 315 5.2 Parallel training setup (multi-task learning)

### 316 5.2.1 Experimental setting

317 The main difference compared to the sequential setup is that multiple dynamics are to be  
 318 learned simultaneously and that task identifiers are passed to the agent for this purpose uti-  
 319 lizing the DynamicsHintObservation wrapper. For the environments walker-walk-v0 and  
 320 HalfCheetah-v4 we train the agent for  $T = 1000000$  and  $T = 2000000$  steps, respectively. The  
 321 steps eventually trained per task  $\mathcal{T}_i$  are evenly distributed. In general we consider only a pair of tasks,  
 322 hence  $\{\mathcal{T}_1, \mathcal{T}_2\}$ .

323 In our multi-task experiments we solely rely on one-hot task identifiers. For example, Sodhani et al.  
 324 [47] communicate tasks from Meta-World [62] to the agent through language descriptions. They  
 325 show that similarities between these descriptions can be exploited for learning. Applying this idea to  
 326 modifications in action effects would yield interesting experiments, which we defer to future research  
 327 endeavors.

### 328 5.2.2 Metrics

329 Again,  $R_t(\pi, \mathcal{T}_i)$  is the episodic return under the policy  $\pi$  on task  $\mathcal{T}_i$  at time step  $t$ . We define a  
 330 performance measure  $\rho_i(t)$  by normalizing the return by its maximum value over all time steps and  
 331 over all tasks. As in the sequential case, Equation 1 applies.



332 **Average performance** Similar to the setup for sequential training, we measure the performance at  
 333 any time step averaged across all tasks. We consider its final performance  $P(t = T)$  for evaluation.  
 334 Equation 2 applies.

335 **Parallel transfer** In the sequential learning setup Chaudhry et al. [6] define an intransigence metric  
 336 for measuring forward transfer by comparing the maximum return for a task trained independently to  
 337 the return achieved while it was trained sequentially [58, 40].

338 For the multi-task learning setup we propose an analogous metric to assess parallel transfer. We  
 339 compare the performance  $\rho_i^b$  for a task trained independently (single-task, baseline) to the performance  
 340  $\rho_i$  achieved while it was trained in parallel to other tasks (multi-task). We assume both, the single-task  
 341 and the multi-task run, train for the same number of steps.  $\rho_i$  and  $\rho_i^b$  are both normalized by the  
 342 overall (single-task or multi-task condition) maximum value.

$$PT = \frac{1}{N} \sum_{i=1}^N PT_i \quad \text{where} \quad PT_i = \rho_i(t = T) - \rho_i^b(t = T) \quad (6)$$

### 343 5.2.3 Pair-wise multi-task

344 We train a SAC agent in parallel on two tasks in the `walker-walk-v0` environment, i.e., on the  
 345 unmodified action dynamic and on a modification of such. We compare the performances of the  
 346 two tasks trained in a multi-task setting to their baseline performances achieved in a single-task  
 347 setting. While we observe hardly any parallel transfer if the `InvertAction` modification is done on  
 348 only one dimension (cf. Figure 4a,  $PT = -0.02$ ), if it is done on all action dimensions the agent  
 349 cannot cope with both tasks in parallel. Figure 4b illustrates an instance of negative parallel transfer  
 350 ( $PT = -0.14$ ). In all our experiments (cf. Table 11 and 12 in the Appendix) we were not able  
 351 to observe positive parallel transfer which might be due to not using a specific multi-task learning  
 352 algorithm.

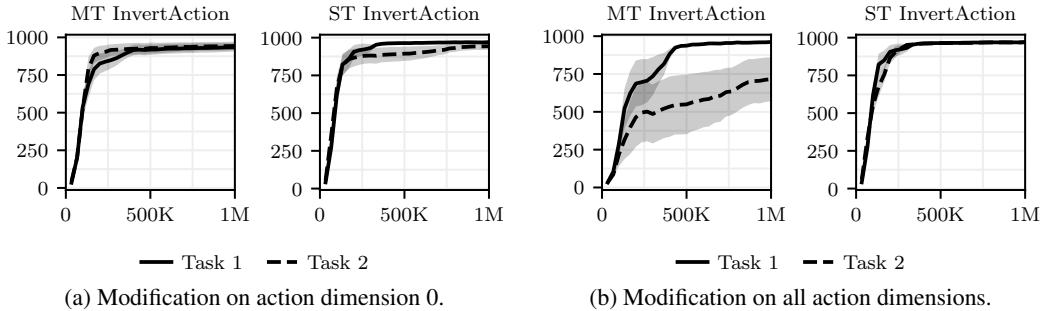


Figure 4: Evaluation returns for a SAC agent trained on the `walker-walk-v0` environment. It independently trains on two tasks in parallel (multi-task, MT) and individually (baseline, single-task, ST). Tasks correspond to the unmodified action dynamic (solid lines) and a `InvertAction` modification (dashed lines).

## 353 6 Conclusions

354 We propose a *Universal Benchmark for Actuation Dynamics Adaptation (UBADA)* comprising a  
 355 set of universal wrappers adhering to the Gymnasium API standard. UBADA can modify arbitrary  
 356 environments and turn them into challenges for both continual (serial) and multi-task (parallel)  
 357 learning scenarios. It focuses specifically on adaptation to changing action dynamics. With our  
 358 experiments we utilize these challenges to advance the understanding and evaluation of RL agents’  
 359 transfer capabilities under continual and sudden dynamic changes. Research on adaptability to  
 360 changing dynamics is crucial for robustness and eventually real-world applications. With our  
 361 benchmark, we place a clear emphasis on dynamic changes associated with variations in action effects,  
 362 and we are confident that UBADA can aid systematic research on that perspective of robustness.

363 **References**

- 364 [1] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,  
365 O. Pieter Abbeel, and W. Zaremba. Hindsight Experience Replay. In *Advances in Neural*  
366 *Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 367 [2] G. Barth-Maron, M. W. Hoffman, D. Budden, W. Dabney, D. Horgan, D. Tb, A. Muldal,  
368 N. Heess, and T. Lillicrap. Distributed Distributional Deterministic Policy Gradients. In  
369 *International Conference on Learning Representations*, 2018.
- 370 [3] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An  
371 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,  
372 2013.
- 373 [4] C. Benjamins, T. Eimer, F. Schubert, A. Mohan, S. Döhler, A. Biedenkapp, B. Rosenhahn,  
374 F. Hutter, and M. Lindauer. Contextualize Me - The Case for Context in Reinforcement Learning.  
375 *Transactions on Machine Learning Research*, 2023.
- 376 [5] A. Carta, L. Pellegrini, A. Cossu, H. Hemati, and V. Lomonaco. Avalanche: A PyTorch library  
377 for deep continual learning. *Journal of Machine Learning Research*, 24(363):1–6, 2023.
- 378 [6] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. S. Torr. Riemannian Walk for Incremental  
379 Learning: Understanding Forgetting and Intransigence. In *Proceedings of the European*  
380 *Conference on Computer Vision (ECCV)*, pages 532–547, 2018.
- 381 [7] M. Chevalier-Boisvert, D. Bahdanau, S. Lahlou, L. Willems, C. Saharia, T. H. Nguyen, and  
382 Y. Bengio. BabyAI: A Platform to Study the Sample Efficiency of Grounded Language Learning,  
383 2019.
- 384 [8] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging Procedural Generation to Bench-  
385 mark Reinforcement Learning. In *Proceedings of the 37th International Conference on Machine*  
386 *Learning*, pages 2048–2056. PMLR, 2020.
- 387 [9] R. de Lazcano, K. Andreas, J. J. Tai, S. R. Lee, and J. Terry. Gymnasium robotics, 2023.
- 388 [10] N. Díaz-Rodríguez, V. Lomonaco, D. Filliat, and D. Maltoni. Don’t forget, there is more than  
389 forgetting: New metrics for Continual Learning. *CoRR*, abs/1810.13166, 2018.
- 390 [11] G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester. An  
391 empirical investigation of the challenges of real-world reinforcement learning, 2021.
- 392 [12] O. Eberhard, J. Hollenstein, C. Pinneri, and G. Martius. Pink Noise Is All You Need: Colored  
393 Noise Exploration in Deep Reinforcement Learning. In *The Eleventh International Conference*  
394 *on Learning Representations*, 2022.
- 395 [13] J. P. Ebert and D. M. Wegner. Time warp: Authorship shapes the perceived timing of actions  
396 and events. *Consciousness and Cognition*, 19(1):481–489, 2010. ISSN 1053-8100. doi:  
397 10.1016/j.concog.2009.10.002.
- 398 [14] B. Ellenberger. PyBullet gymperium, 2018/2019.
- 399 [15] S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic  
400 methods. In *International Conference on Machine Learning*, pages 1582–1591, 2018.
- 401 [16] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft Actor-Critic: Off-Policy Maximum  
402 Entropy Deep Reinforcement Learning with a Stochastic Actor. In *Proceedings of the 35th*  
403 *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- 404 [17] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent  
405 dynamics for planning from pixels. In *International Conference on Machine Learning*, pages  
406 2555–2565, 2019.
- 407 [18] N. Hansen and X. Wang. Generalization in reinforcement learning by soft data augmentation.  
408 In *International Conference on Robotics and Automation*, 2021.
- 409 [19] P. Henderson, W.-D. Chang, F. Shkurti, J. Hansen, D. Meger, and G. Dudek. Benchmark  
410 Environments for Multitask Learning in Continuous Domains, 2017.
- 411 [20] J. Hollenstein, S. Auddy, M. Saveriano, E. Renaudo, and J. Piater. Action noise in off-policy  
412 deep reinforcement learning: Impact on exploration and performance. *Transactions on Machine*  
413 *Learning Research*, 2022. ISSN 2835-8856.

- 414 [21] M. Igl, G. Farquhar, J. Luketina, W. Boehmer, and S. Whiteson. Transient Non-stationarity  
415 and Generalisation in Deep Reinforcement Learning. In *International Conference on Learning*  
416 *Representations*, 2020.
- 417 [22] S. James, Z. Ma, D. Rovick Arrojo, and A. J. Davison. RL Bench: The robot learning benchmark  
418 & learning environment. *IEEE Robotics and Automation Letters*, 2020.
- 419 [23] R. Khraishi and R. Okhrati. Simple Noisy Environment Augmentation for Reinforcement  
420 Learning, 2023.
- 421 [24] R. Kirk, A. Zhang, E. Grefenstette, and T. Rocktäschel. A Survey of Zero-shot Generalisation  
422 in Deep Reinforcement Learning. *Journal of Artificial Intelligence Research*, 76:201–264, 2023.  
423 ISSN 1076-9757. doi: 10.1613/jair.1.14174.
- 424 [25] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan,  
425 J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and  
426 R. Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National*  
427 *Academy of Sciences*, 114(13):3521–3526, 2017. doi: 10.1073/pnas.1611835114.
- 428 [26] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani,  
429 D. Gordon, Y. Zhu, A. Kembhavi, A. Gupta, and A. Farhadi. AI2-THOR: An Interactive 3D  
430 Environment for Visual AI, 2022.
- 431 [27] I. Kostrikov. JAXRL: Implementations of reinforcement learning algorithms in JAX, 2022.
- 432 [28] H. Küttler, N. Nardelli, A. Miller, R. Raileanu, M. Selvatici, E. Grefenstette, and T. Rocktäschel.  
433 The NetHack Learning Environment. In *Advances in Neural Information Processing Systems*,  
434 volume 33, pages 7671–7684. Curran Associates, Inc., 2020.
- 435 [29] E. D. Langlois and T. Everitt. How RL agents behave when their actions are modified. *Pro-*  
436 *ceedings of the AAAI Conference on Artificial Intelligence*, 35(13):11586–11594, 2021. doi:  
437 10.1609/aaai.v35i13.17378.
- 438 [30] J. Leike, M. Martic, V. Krakovna, P. A. Ortega, T. Everitt, A. Lefrancq, L. Orseau, and S. Legg.  
439 AI Safety Gridworlds, 2017.
- 440 [31] M. Liesner, W. Kirsch, and W. Kunde. The interplay of predictive and postdictive components  
441 of experienced selfhood. *Consciousness and Cognition*, 77:102850, 2020. ISSN 1053-8100.  
442 doi: 10.1016/j.concog.2019.102850.
- 443 [32] M. Liesner, W. Kirsch, R. Pfister, and W. Kunde. Spatial action–effect binding depends on  
444 type of action–effect transformation. *Attention, Perception, & Psychophysics*, 82(5):2531–2543,  
445 2020. ISSN 1943-393X. doi: 10.3758/s13414-020-02013-2.
- 446 [33] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.  
447 Continuous control with deep reinforcement learning. In *4th International Conference on*  
448 *Learning Representations, ICLR*, 2016.
- 449 [34] D. Lopez-Paz and M. A. Ranzato. Gradient Episodic Memory for Continual Learning. In  
450 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 451 [35] N. Lucchesi, A. Carta, V. Lomonaco, and D. Bacciu. Avalanche RL: A Continual Reinforcement  
452 Learning Library, 2022.
- 453 [36] M. C. Machado, M. G. Bellemare, E. Talvitie, J. Veness, M. Hausknecht, and M. Bowling.  
454 Revisiting the arcade learning environment: Evaluation protocols and open problems for general  
455 agents. *Journal of Artificial Intelligence Research*, 61(1):523–562, 2018. ISSN 1076-9757.
- 456 [37] J. A. Mendez, M. Hussing, M. Gummadi, and E. Eaton. CompoSuite: A Compositional  
457 Reinforcement Learning Benchmark, 2022.
- 458 [38] J. Parker-Holder, R. Rajan, X. Song, A. Biedenkapp, Y. Miao, T. Eimer, B. Zhang, V. Nguyen,  
459 R. Calandra, A. Faust, F. Hutter, and M. Lindauer. Automated Reinforcement Learning  
460 (AutoRL): A Survey and Open Problems. *Journal of Artificial Intelligence Research*, 74, 2022.  
461 ISSN 1076-9757. doi: 10.1613/jair.1.13596.
- 462 [39] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin,  
463 M. Chociej, P. Welinder, V. Kumar, and W. Zaremba. Multi-Goal Reinforcement Learning:  
464 Challenging Robotics Environments and Request for Research, 2018.

- 465 [40] S. Powers, E. Xing, E. Kolve, R. Mottaghi, and A. Gupta. CORA: Benchmarks, Baselines, and  
466 Metrics as a Platform for Continual Reinforcement Learning Agents. In *Proceedings of The 1st*  
467 *Conference on Lifelong Learning Agents*, pages 705–743. PMLR, 2022.
- 468 [41] A. Raffin. RL baselines3 zoo, 2020.
- 469 [42] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-Baselines3:  
470 Reliable Reinforcement Learning Implementations. *Journal of Machine Learning Research*, 22  
471 (268):1–8, 2021. ISSN 1533-7928.
- 472 [43] D. Rolnick, A. Ahuja, J. Schwarz, T. Lillicrap, and G. Wayne. Experience Replay for Continual  
473 Learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates,  
474 Inc., 2019.
- 475 [44] M. Samvelyan, R. Kirk, V. Kurin, J. Parker-Holder, M. Jiang, E. Hambro, F. Petroni, H. Kuttler,  
476 E. Grefenstette, and T. Rocktäschel. MiniHack the planet: A sandbox for open-ended reinforce-  
477 ment learning research. In *Thirty-Fifth Conference on Neural Information Processing Systems*  
478 *Datasets and Benchmarks Track (Round 1)*, 2021.
- 479 [45] J. Schwarz, W. Czarnecki, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and  
480 R. Hadsell. Progress & Compress: A scalable framework for continual learning. In *Proceedings*  
481 *of the 35th International Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.
- 482 [46] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and  
483 D. Fox. ALFRED: A benchmark for interpreting grounded instructions for everyday tasks. In  
484 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 485 [47] S. Sodhani, A. Zhang, and J. Pineau. Multi-Task Reinforcement Learning with Context-based  
486 Representations, 2021.
- 487 [48] D. Starch. A demonstration of the trial and error method of learning. *Psychological Bulletin*, 7  
488 (1):20–23, 1910. ISSN 1939-1455. doi: 10.1037/h0063796.
- 489 [49] A. Stone, O. Ramirez, K. Konolige, and R. Jonschkowski. The Distracting Control Suite – A  
490 Challenging Benchmark for Reinforcement Learning from Pixels, 2021.
- 491 [50] L. Sun, H. Zhang, W. Xu, and M. Tomizuka. PaCo: Parameter-Compositional Multi-Task  
492 Reinforcement Learning, 2022.
- 493 [51] J. J. Tai, M. Towers, and E. Tower. Shimmy: Gymnasium and PettingZoo wrappers for  
494 commonly used environments. Zenodo, 2023.
- 495 [52] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *2012*  
496 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.  
497 doi: 10.1109/IROS.2012.6386109.
- 498 [53] M. Towers, J. K. Terry, A. Kwiatkowski, J. U. Balis, G. de Cola, T. Deleu, M. Goulão,  
499 A. Kallinteris, A. KG, M. Krimmel, R. Perez-Vicente, A. Pierré, S. Schulhoff, J. J. Tai, A. T. J.  
500 Shen, and O. G. Younis. Gymnasium, 2023.
- 501 [54] S. Tunyasuvunakool, A. Muldal, Y. Doron, S. Liu, S. Bohez, J. Merel, T. Erez, T. Lillicrap,  
502 N. Heess, and Y. Tassa. Dm<sub>c</sub>ontrol: Software and tasks for continuous control. *Software*  
503 *Impacts*, 6:100022, 2020. ISSN 2665-9638. doi: 10.1016/j.simpa.2020.100022.
- 504 [55] G. E. Uhlenbeck and L. S. Ornstein. On the theory of the brownian motion. *Physical Review*,  
505 36(5):823–841, 1930. doi: 10.1103/PhysRev.36.823.
- 506 [56] H. van Seijen, H. Nekoei, E. Racah, and S. Chandar. The LoCA regret: A consistent metric  
507 to evaluate model-based behavior in reinforcement learning. In *Proceedings of the 34th*  
508 *International Conference on Neural Information Processing Systems*, 2020.
- 509 [57] P. Wawrzyński. A Cat-Like Robot Real-Time Learning to Run. In *Adaptive and Natural Comput-*  
510 *ing Algorithms*, volume 5495, pages 380–390. Springer Berlin Heidelberg, Berlin, Heidelberg,  
511 2009. ISBN 978-3-642-04920-0 978-3-642-04921-7. doi: 10.1007/978-3-642-04921-7\_39.
- 512 [58] M. Wołczyk, M. Zając, R. Pascanu, Ł. Kuciński, and P. Miłoś. Continual World: A Robotic  
513 Benchmark For Continual Reinforcement Learning, 2021.
- 514 [59] D. Yarats, A. Zhang, I. Kostrikov, B. Amos, J. Pineau, and R. Fergus. Improving sample  
515 efficiency in model-free reinforcement learning from images. 2019.

- 516 [60] D. Yarats, I. Kostrikov, and R. Fergus. Image augmentation is all you need: Regularizing deep  
517 reinforcement learning from pixels. In *International Conference on Learning Representations*,  
518 2021.
- 519 [61] T. Yu, S. Kumar, A. Gupta, S. Levine, K. Hausman, and C. Finn. Gradient Surgery for  
520 Multi-Task Learning, 2020.
- 521 [62] T. Yu, D. Quillen, Z. He, R. Julian, A. Narayan, H. Shively, A. Bellathur, K. Hausman,  
522 C. Finn, and S. Levine. Meta-World: A Benchmark and Evaluation for Multi-Task and Meta  
523 Reinforcement Learning, 2021.
- 524 [63] L. Zintgraf, S. Schulze, C. Lu, L. Feng, M. Igl, K. Shiarlis, Y. Gal, K. Hofmann, and S. Whiteson.  
525 VariBAD: Variational Bayes-Adaptive Deep RL via Meta-Learning. *Journal of Machine*  
526 *Learning Research*, 22(289):1–39, 2021. ISSN 1533-7928.

## 527 **A More examples**

528 A well known example for switched action dynamics is the backward bicycle, which refers to a  
529 modified bicycle that has its handlebars connected to the front wheel over a gear, so that turning the  
530 handlebars to the right makes the wheel turn to the left, and vice versa. This modification creates  
531 a counterintuitive steering mechanism where the bike responds opposite to the rider’s expectations.  
532 Riding such a bike can be a challenging and it highlights how difficult it can be for individuals to  
533 adapt to new ways of thinking or doing things, even when presented with clear evidence of the need  
534 for change.

535 Note the difference to the mirror drawing or tracing task, a famous test in psychology for which  
536 participants have to draw or trace a figure (such as a shape or a pattern) on the drawing surface by  
537 looking only at the reflection of their hand in the mirror. The catch is that the visual-motor feedback is  
538 reversed due to the mirror, meaning that if the participant moves their hand to the right, the reflection  
539 in the mirror appears to move to the left on the drawing surface [48]. Though being closely related,  
540 in that case the action effect remains the same, instead only the observation is reversed.

## 541 **B Extended background**

### 542 **B.1 Transfer Learning in RL**

543 Transfer learning in the context of RL involves leveraging knowledge gained from one task to improve  
544 the learning or performance of a related but different task. This is especially relevant when the tasks  
545 share some underlying structure or features.

546 On a basic level one can distinguish two kinds of transfer: Positive transfer, where the knowledge  
547 gained on a source task enhances the learning or performance on the target task and negative transfer,  
548 where the transferred knowledge hinders performance on the target task. Obviously positive transfer  
549 is desired whilst negative transfer is a potential risk in transfer learning, and careful design is needed  
550 to mitigate it.

551 For our purpose it may also be beneficial to distinguish parallel and serial transfer. Parallel transfer  
552 implies the simultaneous transfer of knowledge or skills across multiple tasks. In an RL context, this  
553 might involve training on several tasks concurrently, with the expectation that knowledge gained  
554 in one task can benefit learning in others. This is the intention behind multi-task RL. The subtle  
555 difference in serial transfer is that future tasks are unknown, thus it refers to a sequential transfer of  
556 knowledge or skills. In the broad field of meta RL, this corresponds to transferring knowledge from a  
557 source task to an unknown target task. In the context of continual RL, serial transfer translates to  
558 the learning of a series of tasks one after the other. The scheduling of tasks is further specified in  
559 curriculum learning, where each task builds on the knowledge gained from previous ones, possibly  
560 with increasing difficulty or complexity.

561 In the serial transfer learning domain, specifically for the continual learning case one may further  
562 make a difference between forward and backward transfer. Forward transfer occurs when knowledge  
563 or skills acquired from learning on earlier tasks positively influence the learning or performance on  
564 future tasks. Backward transfer, on the other hand, occurs when knowledge or skills acquired from  
565 learning on later tasks positively influence the learning or performance on earlier tasks.

### 566 **B.2 Multi-task RL**

567 The objective of multi-task RL is to acquire a unified policy, denoted as  $\pi(a|s, z)$ , where  $z$  represents  
568 the task ID or an encoding thereof. This policy is designed to maximize the average expected return  
569 across all considered tasks. Task information can be communicated to the policy through various  
570 means, such as language or a one-hot task identification encoding, which is supplied in addition to  
571 the current state.

572 One seeks to understand how effectively the policy generalizes and performs across a spectrum of  
573 related tasks. Multi-task RL algorithms attempt to make use of shared knowledge and skills across  
574 these tasks, enhancing the efficiency of the learning process, i.e., learn a set of tasks more quickly  
575 and more proficiently than learning them independently. This can be considered a parallel transfer of  
576 knowledge. The constant access to all tasks is equivalent to ignoring non-stationarity, while continual

577 RL focuses on just that. Multi-task RL algorithms are commonly evaluated by considering their  
 578 average performance across all training tasks, as opposed to meta RL, which utilizes separate test  
 579 tasks for assessment.

### 580 B.3 Continual RL

581 Continual RL is a field of study dedicated to creating algorithms that can effectively adapt to changing  
 582 environments, i.e., they are capable of handling non-stationarity. The goal is to develop agents that  
 583 can progressively learn new skills and tackle unfamiliar tasks without neglecting what they have  
 584 learned before. This ability to adapt continuously over extended periods is often referred to as lifelong  
 585 learning or endless adaptation.

586 In the training process, a series of tasks is presented to the system. The transitions between these  
 587 tasks may be seamless and unknown to the agent. When evaluating these systems, researchers  
 588 are often interested in assessing the agent’s ability to apply previously acquired knowledge to new  
 589 tasks (forward transfer), as well as retaining knowledge when faced with new challenges (backward  
 590 transfer), all while avoiding catastrophic forgetting. Forward and backward transfer can be considered  
 591 as some kind of serial transfer as opposed to parallel transfer as in the multi-task setting.

## 592 C Action wrapper overview

Table 1: An overview of action effect wrappers including a brief description and motivation.

Wrapper	Description	Motivation
InvertAction	Inverts one/all dimension(s) of the continuous action.	backwards bicycle
ScaleAction	Scales one/all dimension(s) of the continuous action.	sim2real gap
OffsetAction	Adds an offset to one/all dimension(s) of the continuous action.	sim2real gap, systematic error
NoiseAction	Adds gaussian noise to one/all dimension(s) of the continuous action.	sim2real gap, random error
SineNoiseAction	Adds gaussian noise overlaid with a sine offset to one/all dimension(s) of the continuous action.	sim2real gap, random error with systematic component
ZeroAction	Randomly zeros one/all dimension(s) of the continuous action for a defined number of steps.	loose contact
RepeatAction	Randomly repeats one/all dimension(s) of the continuous/discrete action for a defined number of steps.	loose contact
SwapAction	Swaps one dimension of the continuous/discrete action with another (randomly picked) or shuffles all dimensions. Changed order of dimensions is kept.	interchanged cables

## 593 D Hyperparameters

594 We use implementations provided by Kostrikov [27] (JAXRL2) and Raffin et al. [42] (Stable-Baselines  
 595 3). We intended not to tune the hyperparameter and rather kept them the default values. An overview  
 596 is provided in Table 2.

597 Additional remarks: For the `walker-walk-v0(pixel)` environment typically an action repeat value  
 598 of 2 is used [17]. In DrQ, again for `walker-walk-v0(pixel)`, we use the standard choice in

599 JAXRL2 for the image encoder which is originated from D4PG [2]. For the reduction of the twinned  
600 critic we use the mean Q-value in DrQ and the minimum Q-value in SAC which is both the default  
601 configuration. For the goal-conditioned environments `FetchReach-v2` and `FetchPush-v2` we use  
602 SAC as the learner and use HER [1] to sample from the replay buffer. For HER we set the number of  
603 additional, virtual goals sampled per real goal to 4 which is again the default in Stable-Baselines 3.  
604 For the hidden dimensions of actor and critic we used feasible values in accordance to Raffin [41].

Table 2: Hyperparameters for the SAC base used in the experiments.

Parameter	Value
Buffer capacity	1000000
Batch size	256
Discount $\gamma$	0.99
Optimizer	Adam
Critic LR	0.0003
Actor LR	0.0003
Temperature LR	0.0003
Critic soft target update $\tau$	0.005
Init temperature (SAC)	1.0
Init temperature (DrQ)	0.1
Hidden dims (general)	(256, 256)
Hidden dims ( <code>FetchReach-v2</code> )	(64, 64)
Hidden dims ( <code>FetchPush-v2</code> )	(256, 256, 256)

## 605 E Classification of the modifications

606 The purpose of the experiments in Section 5 and G is to illustrate the universal applicability of  
607 action modifications across varied environments, hence they are applicable for different problem  
608 statements. In some cases, the wrappers yield fundamentally different modifications. `NoiseAction`  
609 and `OffsetAction` might be considered comparable in their functionality, whereas `InvertAction`  
610 enables the investigation of a fundamentally different problem. They are intended to reveal the  
611 limits of the agent, in particular with regard to dimensions of robustness and transfer capabilities.  
612 It is not the sole aim to simulate realistic scenarios. For example it is hard to imagine a real-world  
613 scenario where the `SwapAction` wrapper seems appropriate. In Table 1 in the Appendix we provide  
614 the example of interchanged cables which is intended to elucidate the functionality but is arguably  
615 far-fetched. Still, in an idealized world it makes for an interesting case to investigate how the agent  
616 behaves under these circumstances.

617 Based on the experimental results in Section 5 and Section G one may distinguish two main groups.  
618 First, `SwapAction` and `InvertAction`, although the modifications seem minimal, hardly any transfer  
619 is possible for the baseline agent. In the sequential setup it basically has to relearn from scratch (cf.  
620 Figure 5), possibly the return drops to the initial level (cf. Figure 2, respectively on the right). For  
621 the other modifications which may form the second group, namely `ScaleAction`, `OffsetAction`,  
622 `NoiseAction`, `SineNoiseAction`, `ZeroAction`, `RepeatAction`, transfer seems natural. Though,  
623 the experiments show that also with modifications from this group the adaptation can become arbitrary  
624 hard, depending on the configuration value and whether the modification is applied on only one  
625 dimension or all.

## 626 F Relearning in a new dynamic

627 Analogously to the one dynamic switch experiment in Section 5.1.3 we trained a SAC agent for  $T_{\mathcal{T}} =$   
628 20000 steps to succeed on the `FetchReach-v2` environment (cf. Figure 5a). In that environment  
629 the agent is supposed to reach a sampled goal with a robotic arm. Then we modified the action



630 dynamic using the `InvertAction` wrapper such that the action dimension which corresponds to  
 631 the x coordinate is inverted. The agent trains for another  $T_{\mathcal{T}} = 20000$  steps on the new dynamic.  
 632 Figure 5b shows trajectories of the goal reaching robot arm before (stage 0) and after (stages 1-5)  
 633 the dynamic switch. As expected, right after the switch (no relearning, stage 1) the trajectories are  
 634 mirrored in the x coordinate, while the position on the y coordinate is intact. The agent seems to  
 635 relearn the behavior corresponding to the x coordinate from scratch, respectively, it has to compensate  
 636 for the mirroring which seems to gradually improve over the stages 2-5.

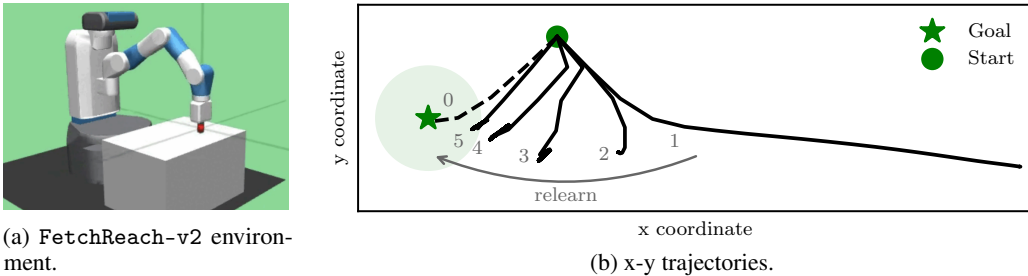


Figure 5: Relearning behavior of a SAC agent in the FetchReach-v2 environment. It is first trained on the default dynamic and then a switch is established using the `InvertAction` wrapper on the dimension 0 which corresponds to the x coordinate. Trajectories (in x-y coordinates) of the goal reaching robot arm before (stage 0, dashed trajectory) and after (stages 1-5, solid trajectory) the dynamic switch are shown.

## 637 G Extended results for experiments

638 In the Figures and Tables within this section we show extended results for a variety of task combina-  
 639 tions, within the sequential and the multi-task setup. We present curves for training and evaluation  
 640 returns. While the training curve only provides the overall view, evaluation returns for individual  
 641 tasks allow for a more complete analysis, it allows for calculating the defined metrics in Section 5.

642 We consider five environments or variants of such: The first two are the goal-conditioned environ-  
 643 ments FetchReach-v2 and FetchPush-v2 initially developed by Plappert et al. [39] and currently  
 644 maintained by de Lazcano et al. [9]. The third is HalfCheetah-v4 based on work by Wawrzyński  
 645 [57] and available through Towers et al. [53]. The fourth and fifth are walker-walk-v0 from DM  
 646 Control [54], utilized via Towers et al. [53] and Tai et al. [51], with, respectively, proprioceptive and  
 647 visual observation inputs. We denote the latter as walker-walk-v0(pixel).

648 In the sequential training setup, for the environments FetchReach-v2, FetchPush-v2,  
 649 walker-walk-v0, walker-walk-v0(pixel) and HalfCheetah-v4 we train each task  $\mathcal{T}_i$  for  
 650  $T_{\mathcal{T}} = 100000$ ,  $T_{\mathcal{T}} = 500000$ ,  $T_{\mathcal{T}} = 500000$ ,  $T_{\mathcal{T}} = 500000$  or  $T_{\mathcal{T}} = 1000000$  steps, respec-  
 651 tively. The choice of  $T_{\mathcal{T}}$  is made to ensure that, at the very least, the unmodified environment can be  
 652 trained to a reasonable degree. We only consider sequences of tasks where  $T_{\mathcal{T}}$  is constant for each  
 653 individual task.

654 In the multi-task training setup, for the environments walker-walk-v0 and HalfCheetah-v4 we  
 655 train the agent for  $T = 1000000$  and  $T = 2000000$  steps, respectively. The steps eventually trained  
 656 per task  $\mathcal{T}_i$  are evenly distributed. In general we consider only a pair of tasks, hence  $\{\mathcal{T}_1, \mathcal{T}_2\}$ .

657 An overview for quicker access to all experiment variants is provided in Table 3.

Table 3: An overview of the extended experiments with references to the corresponding performance figures and metric tables.

Setting	Environment	Algorithm	Returns	Metrics
Sequential, one switch, dimension 0	walker-walk-v0	SAC	Figure 6	Table 4
Sequential, one switch, all dimensions	walker-walk-v0	SAC	Figure 7	Table 4
Sequential, one switch, dimension 0	walker-walk-v0 (pixel)	DrQ	Figure 8	Table 5
Sequential, one switch, all dimensions	walker-walk-v0 (pixel)	DrQ	Figure 9	Table 5
Sequential, one switch, dimension 0	HalfCheetah-v4	SAC	Figure 10	Table 6
Sequential, one switch, all dimensions	HalfCheetah-v4	SAC	Figure 11	Table 6
Sequential, one switch, dimension 0	FetchReach-v2	SAC+HER	Figure 12	Table 7
Sequential, one switch, all dimensions	FetchReach-v2	SAC+HER	Figure 13	Table 7
Sequential, one switch, dimension 0	FetchPush-v2	SAC+HER	Figure 14	Table 8
Sequential, one switch, all dimensions	FetchPush-v2	SAC+HER	Figure 15	Table 8
Sequential, one parallel switch, combine ScaleAction, all dimensions	walker-walk-v0	SAC	Figure 16	-
Sequential, one parallel switch, combine OffsetAction, all dimensions	walker-walk-v0	SAC	Figure 17	-
Sequential, continual adaptation, 10 tasks, dimension 0	walker-walk-v0	SAC	Figure 18	Table 9
Sequential, continual adaptation, 10 tasks, all dimensions	walker-walk-v0	SAC	Figure 19	Table 9
Sequential, continual adaptation, 10 tasks, dimension 0	HalfCheetah-v4	SAC	Figure 20	Table 10
Sequential, continual adaptation, 10 tasks, all dimensions	HalfCheetah-v4	SAC	Figure 21	Table 10
Sequential, continual adaptation, 500 tasks, dimension 0	walker-walk-v0	SAC	Figure 22	-
Sequential, continual adaptation, 500 tasks, all dimensions	walker-walk-v0	SAC	Figure 23	-
Multi-task, dimension 0	walker-walk-v0	SAC	Figure 24	Table 11
Multi-task, all dimensions	walker-walk-v0	SAC	Figure 25	Table 11
Multi-task, dimension 0	HalfCheetah-v4	SAC	Figure 26	Table 12
Multi-task, all dimensions	HalfCheetah-v4	SAC	Figure 27	Table 12

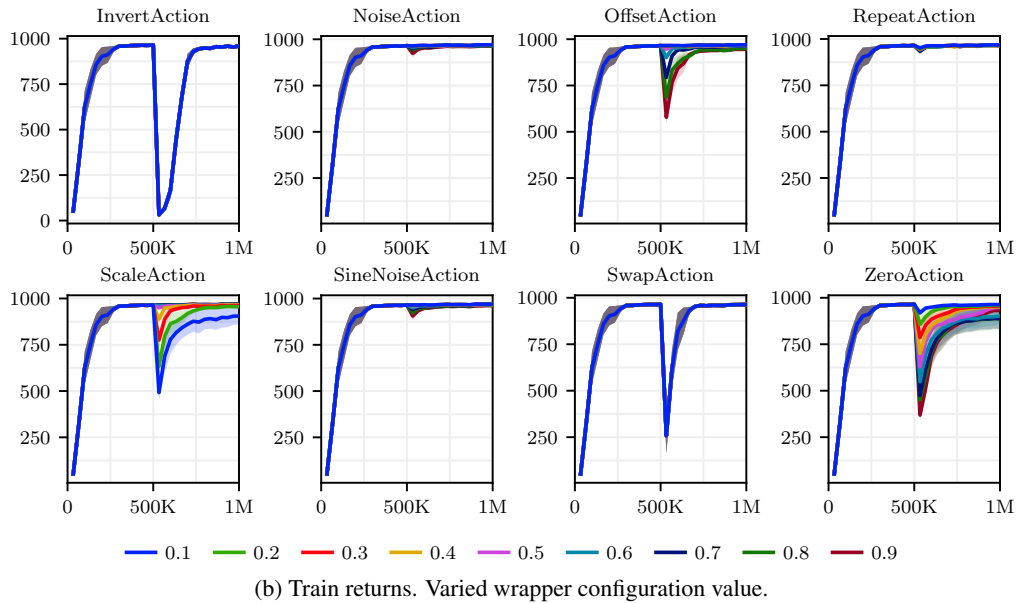
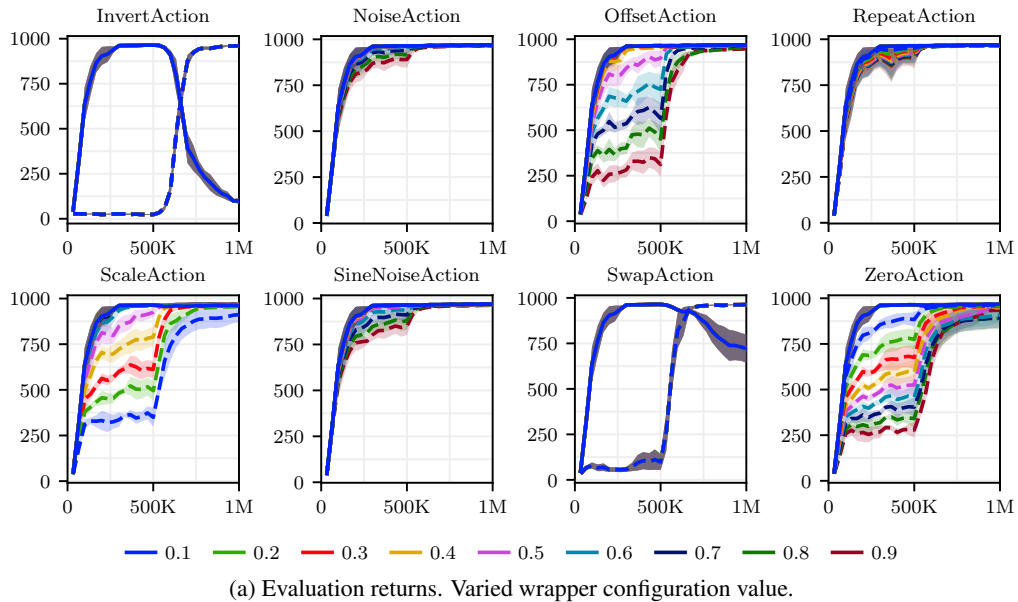


Figure 6: Returns for the **sequential** training setup. A SAC agent is trained on the `walker-walk-v0` environment. **One dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.

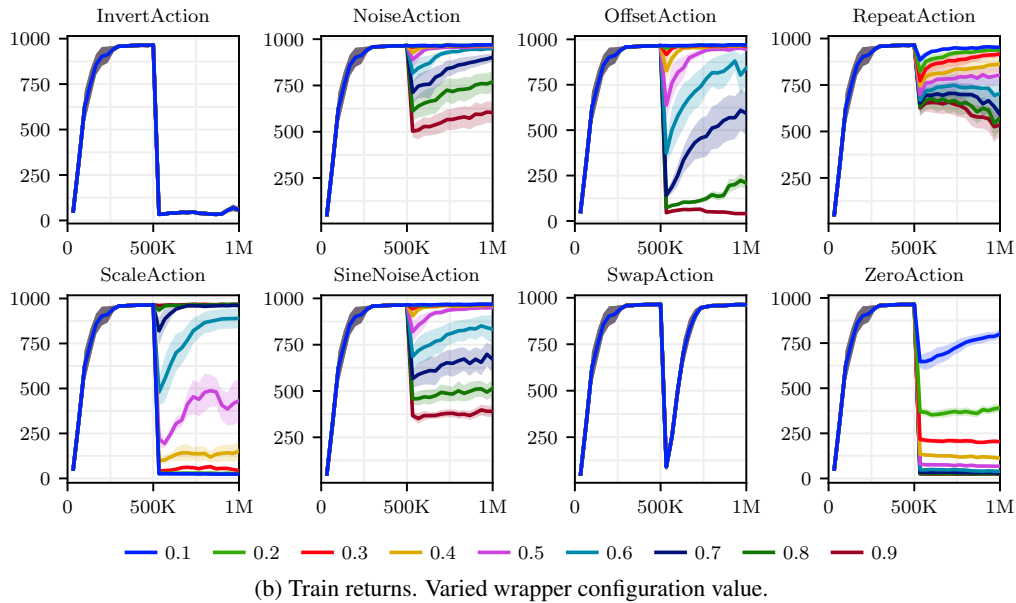
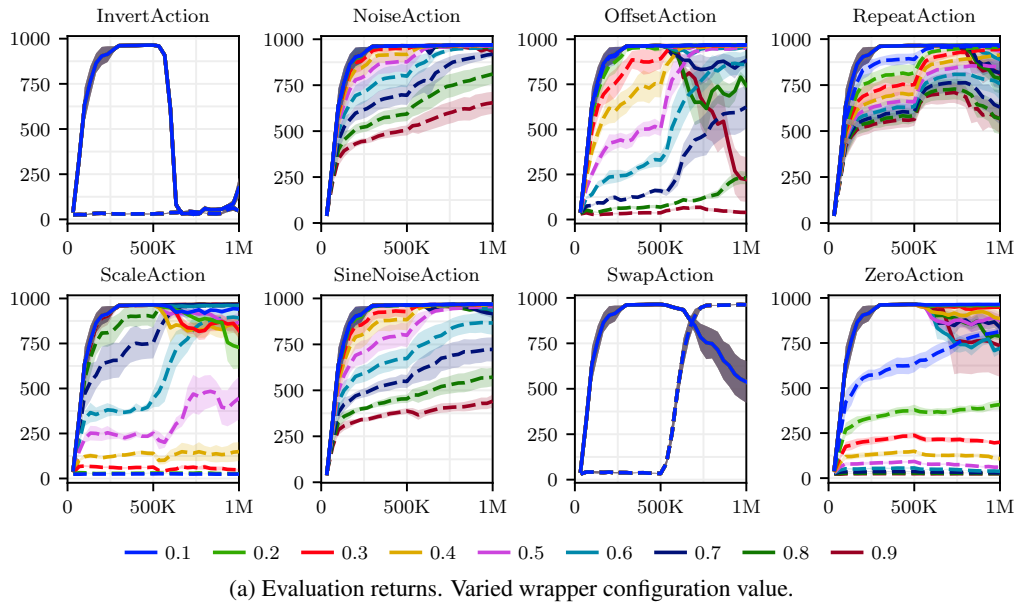
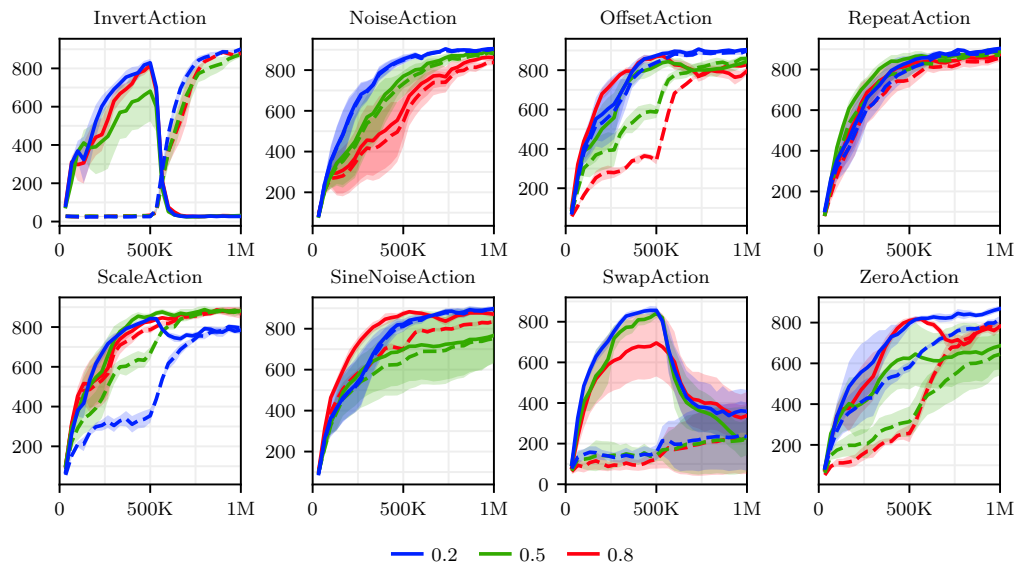
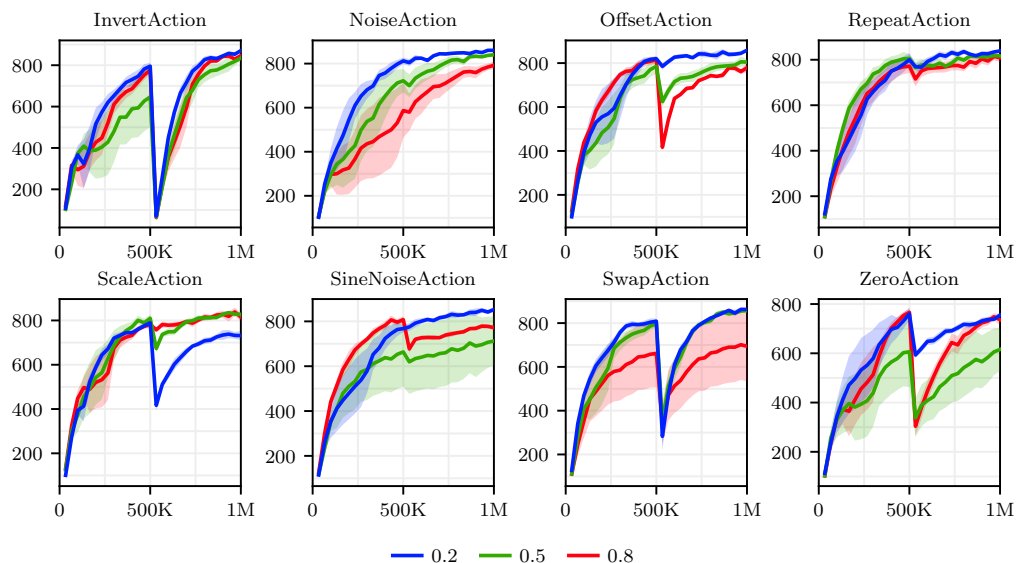


Figure 7: Returns for the **sequential** training setup. A SAC agent is trained on the `walker-walk-v0` environment. **One dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

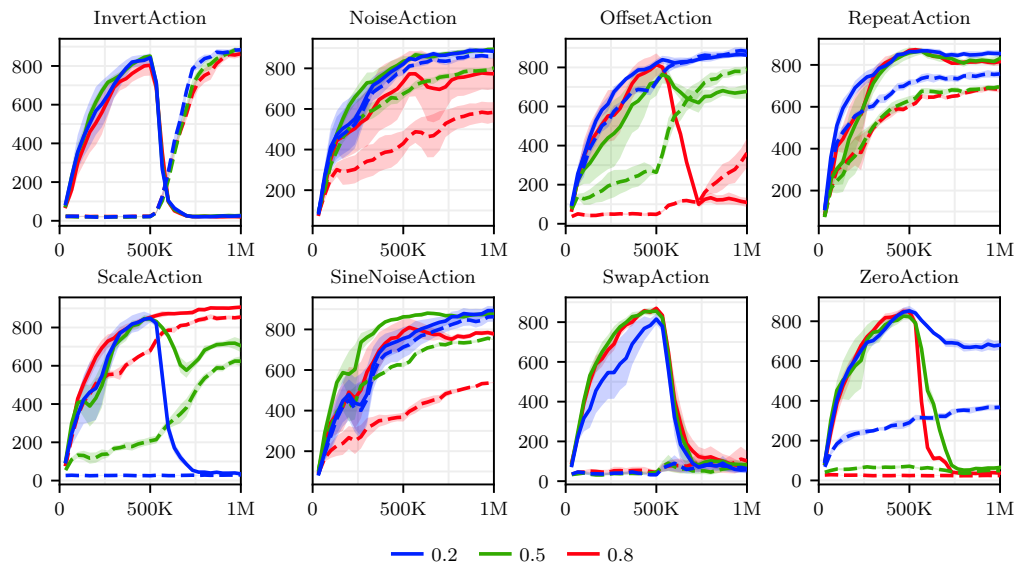


(a) Evaluation returns. Varied wrapper configuration value.

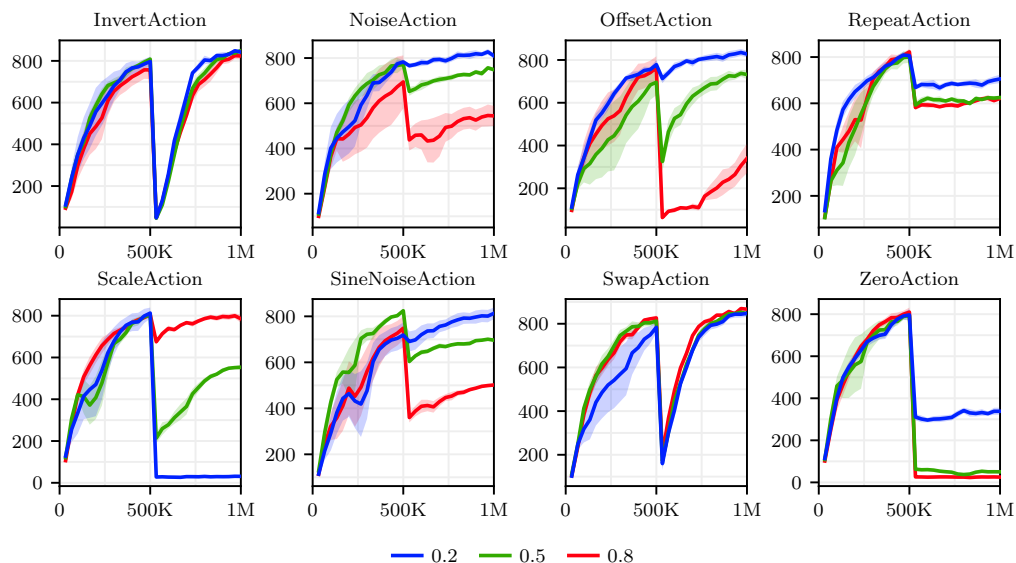


(b) Train returns. Varied wrapper configuration value.

Figure 8: Returns for the **sequential** training setup. A **DrQ** agent is trained on the `walker-walk-v0(pixel)` environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.



(a) Evaluation returns. Varied wrapper configuration value.



(b) Train returns. Varied wrapper configuration value.

Figure 9: Returns for the **sequential** training setup. A **DrQ** agent is trained on the `walker-walk-v0(pixel)` environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

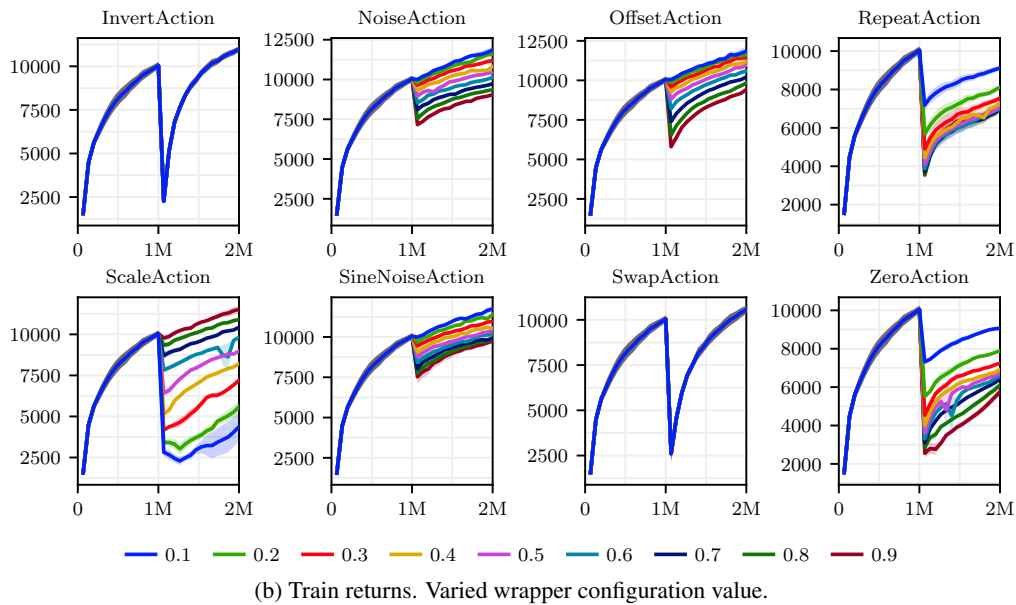
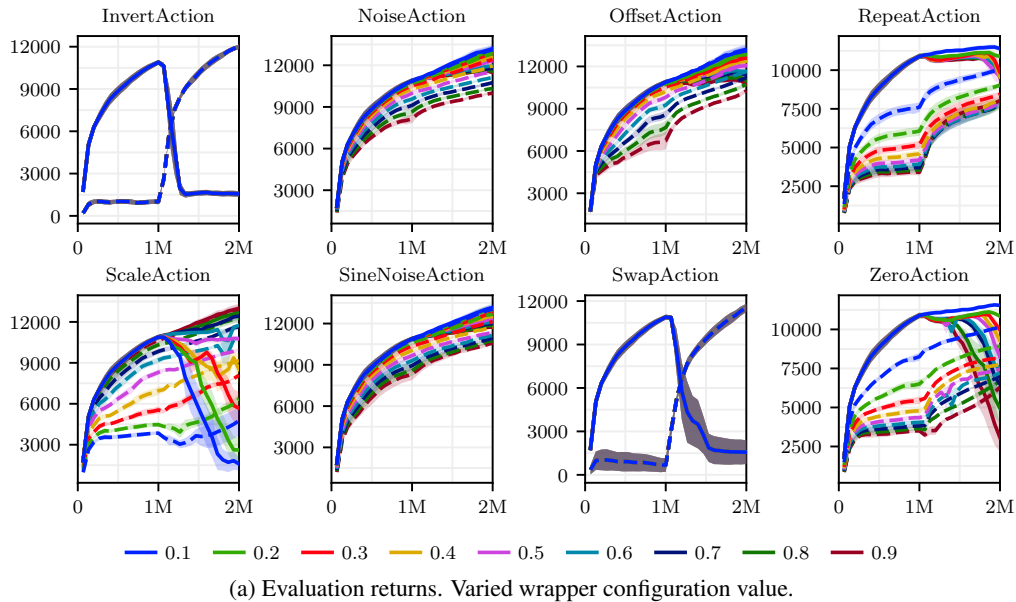


Figure 10: Returns for the **sequential** training setup. A SAC agent is trained on the `HalfCheetah-v4` environment. **One dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.

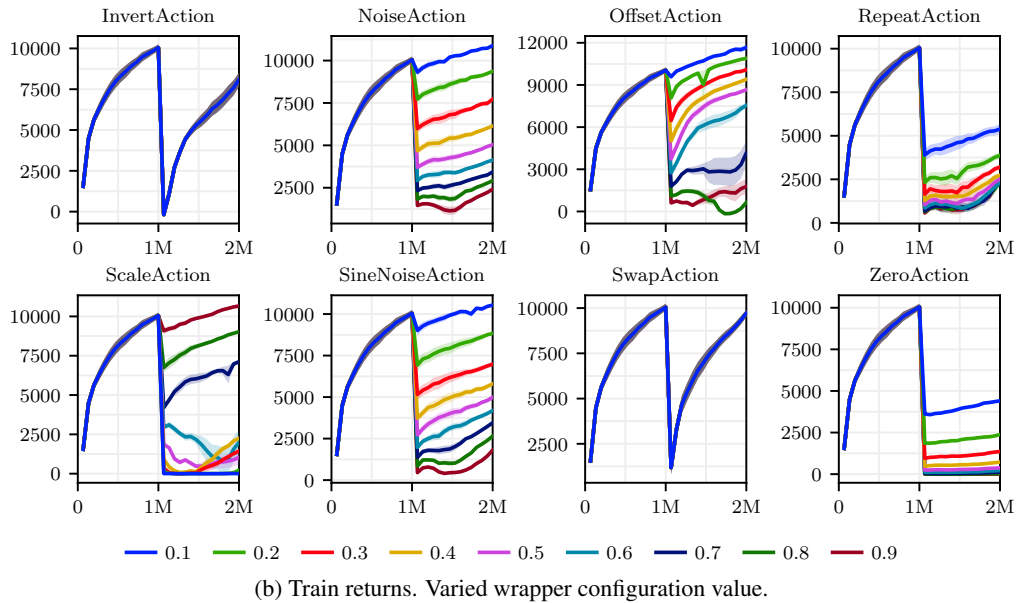
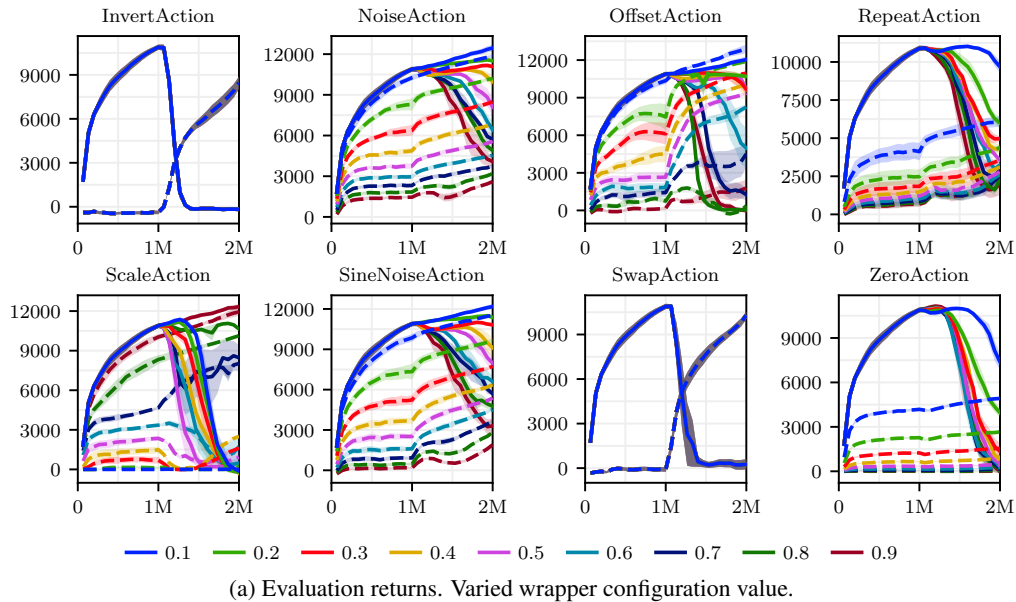


Figure 11: Returns for the **sequential** training setup. A SAC agent is trained on the `HalfCheetah-v4` environment. **One dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.



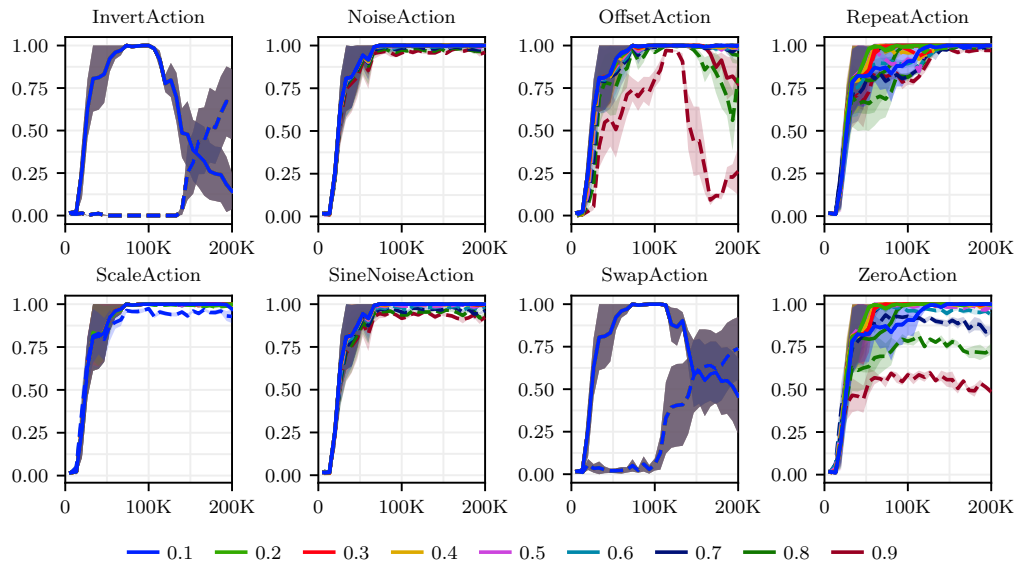


Figure 12: Evaluation success rates for the **sequential** training setup. A SAC agent is trained with **HER** on the goal-conditioned FetchReach-v2 environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.

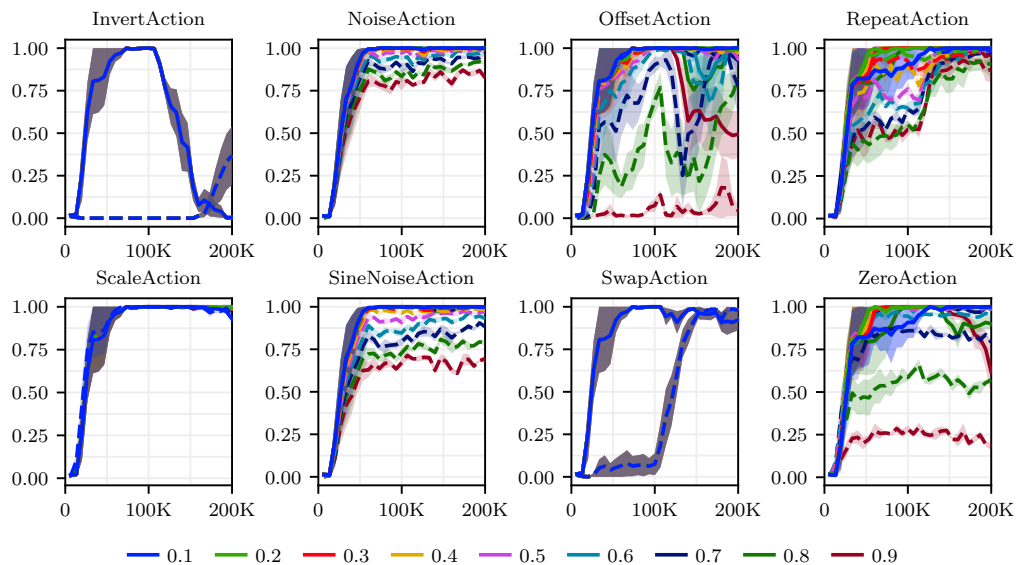


Figure 13: Evaluation success rates for the **sequential** training setup. A SAC agent is trained with **HER** on the goal-conditioned FetchReach-v2 environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

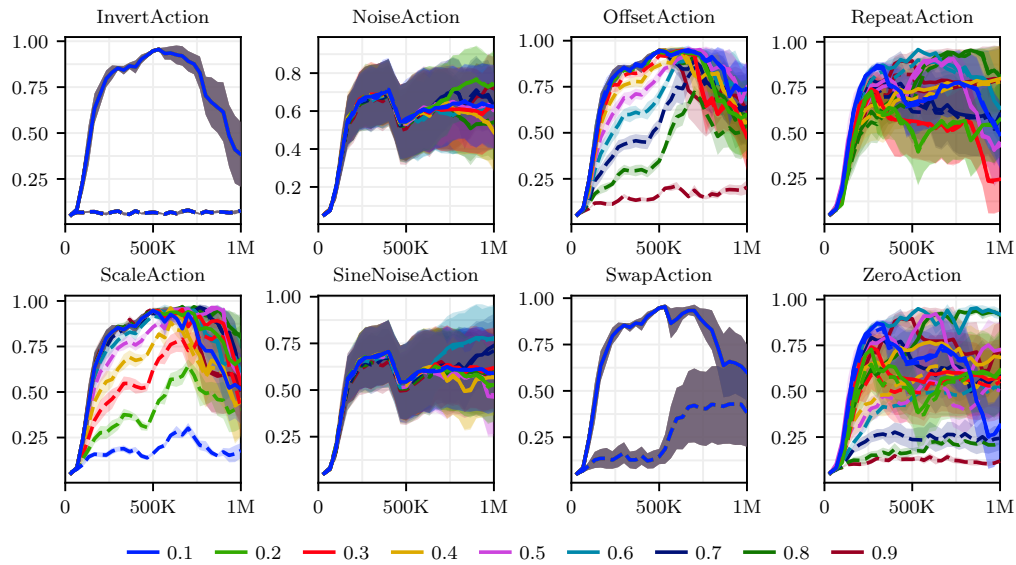


Figure 14: Evaluation success rates for the **sequential** training setup. A SAC agent is trained with **HER** on the goal-conditioned FetchPush-v2 environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.

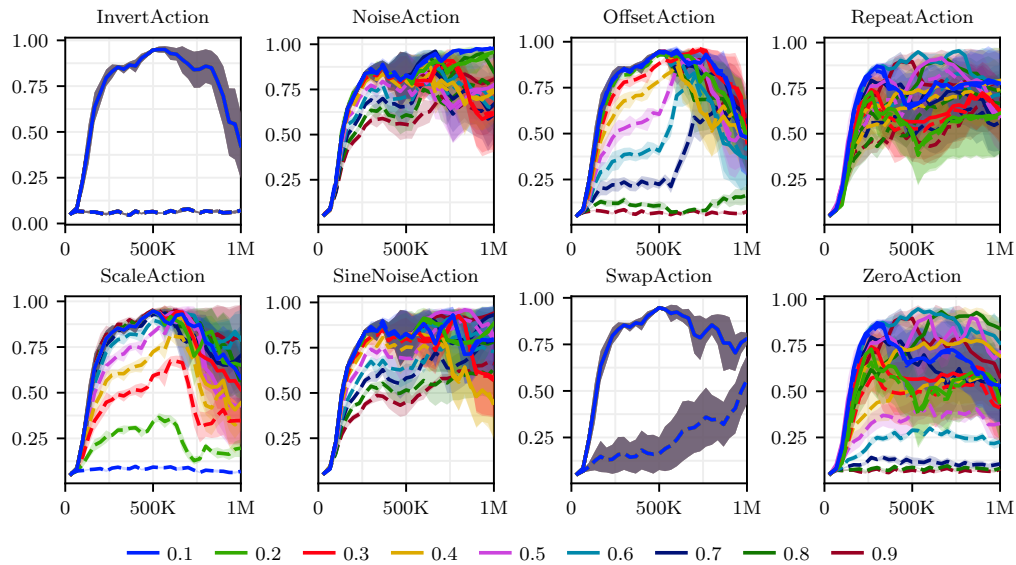


Figure 15: Evaluation success rates for the **sequential** training setup. A SAC agent is trained with **HER** on the goal-conditioned FetchPush-v2 environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

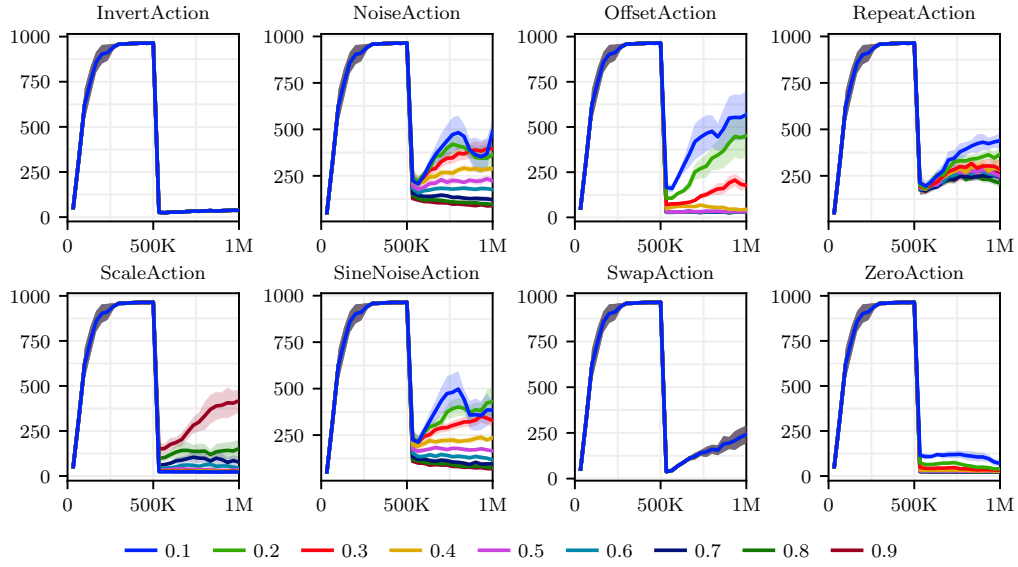
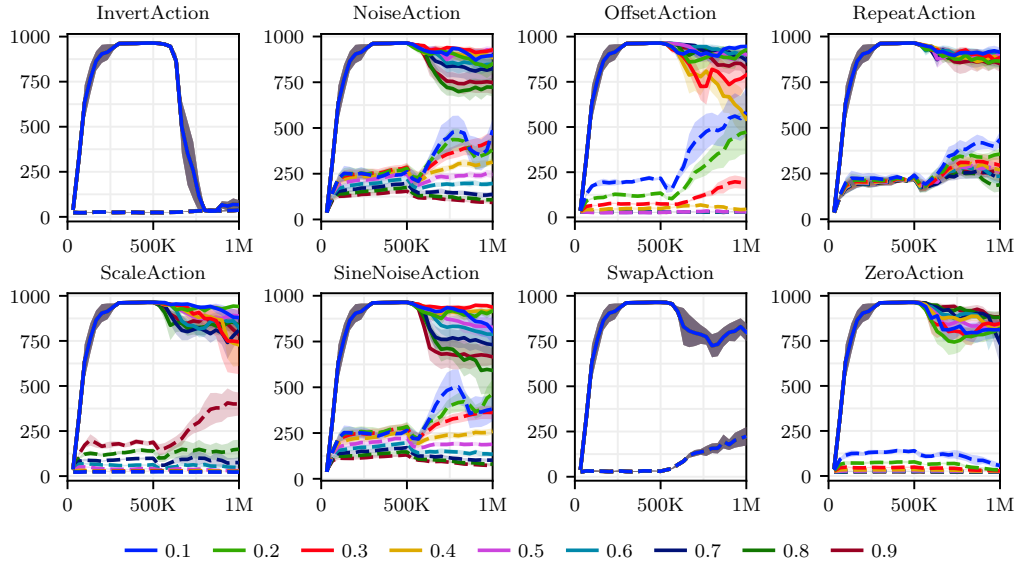


Figure 16: Returns for the **sequential** training setup. A SAC agent is trained on the walker-walk-v0 environment. **One parallel dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. Two parallel modifications are applied. One modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. On top of that, the ScaleAction wrapper is applied with a scaling value of 0.5 representing a combination of changes. The modifications are done on **all action dimensions**.

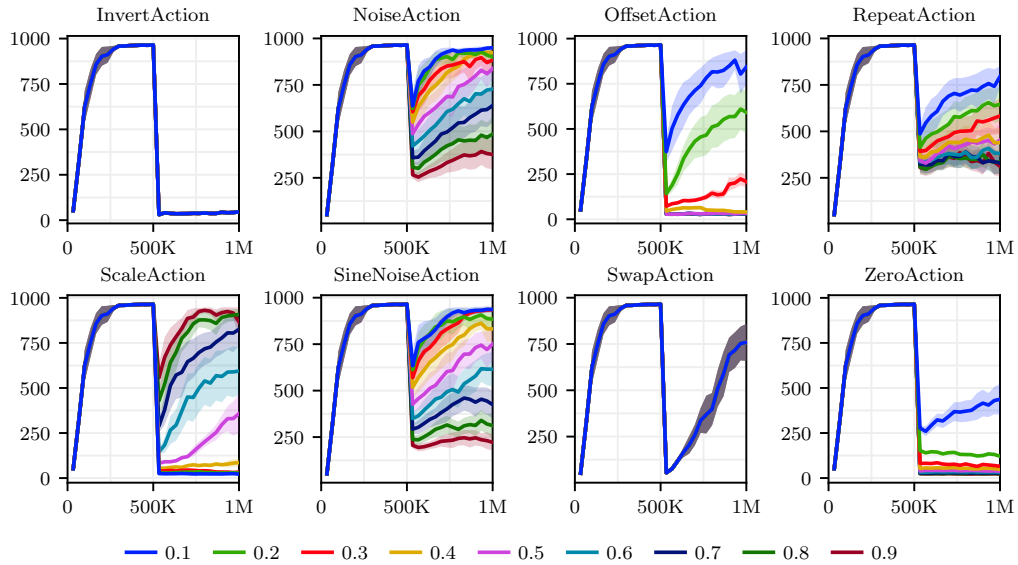
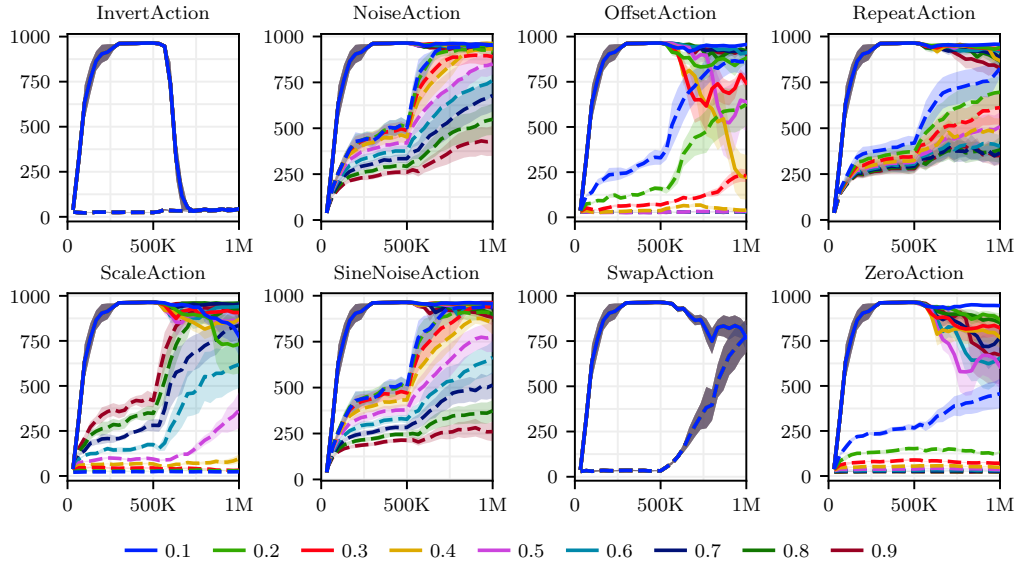
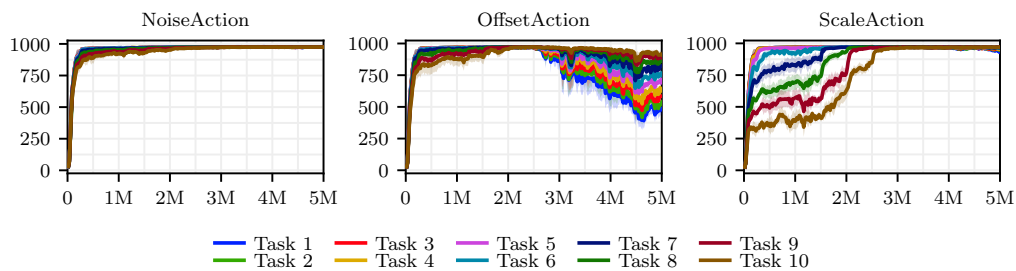
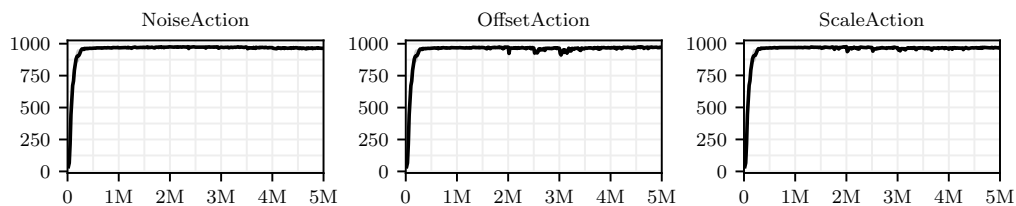


Figure 17: Returns for the **sequential** training setup. A SAC agent is trained on the `walker-walk-v0` environment. **One parallel dynamics switch:** The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. Two parallel modifications are applied. One modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. On top of that, the `OffsetAction` wrapper is applied with an offset value of 0.5 representing a combination of changes. The modifications are done on **all action dimensions**.

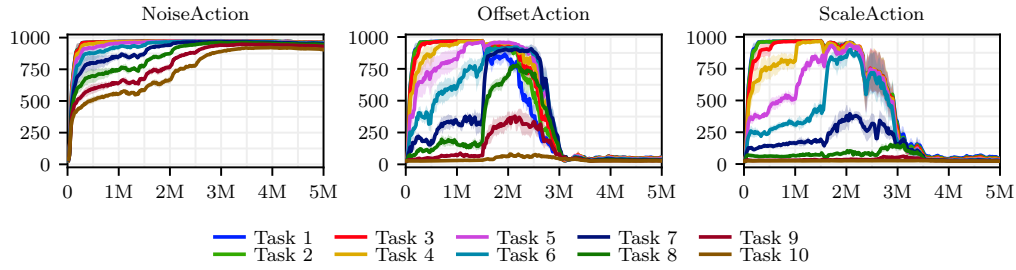


(a) Evaluation returns. Different tasks based on varied wrapper configuration values.

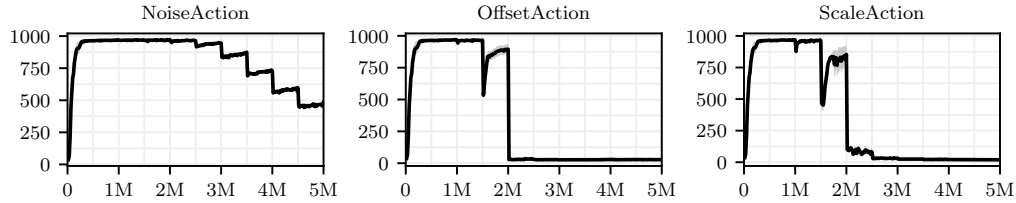


(b) Train returns.

Figure 18: Returns for the **sequential** training setup. A **SAC** agent is trained on the `walker-walk-v0` environment. **Continual adaptation:** The agent trains for  $T/10$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{10}\}$ . For `NoiseAction` tasks are determined by an action wrapper value in  $\{0.1, 0.2, \dots, 0.9\}$ , for `OffsetAction` in  $\{-0.1, -0.2, \dots, -0.9\}$ , for `ScaleAction` in  $\{0.9, 0.8, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. The modification is done on **action dimension 0**.

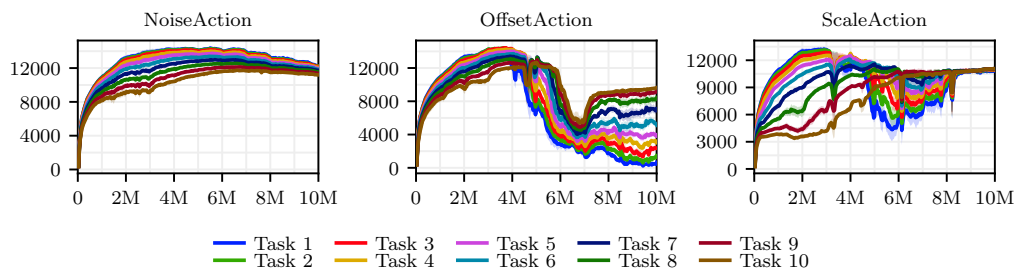


(a) Evaluation returns. Different tasks based on varied wrapper configuration values.

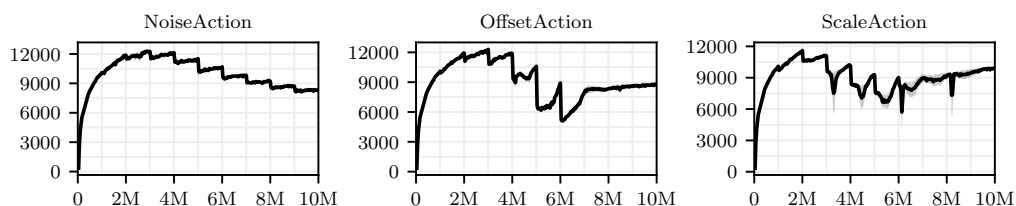


(b) Train returns.

Figure 19: Returns for the **sequential** training setup. A **SAC** agent is trained on the `walker-walk-v0` environment. **Continual adaptation:** The agent trains for  $T/10$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{10}\}$ . For `NoiseAction` tasks are determined by an action wrapper value in  $\{0.1, 0.2, \dots, 0.9\}$ , for `OffsetAction` in  $\{-0.1, -0.2, \dots, -0.9\}$ , for `ScaleAction` in  $\{0.9, 0.8, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. The modification is done on **all action dimensions**.

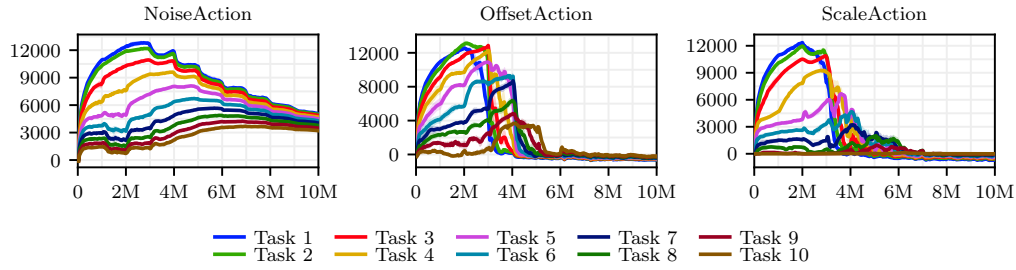


(a) Evaluation returns. Different tasks based on varied wrapper configuration values.

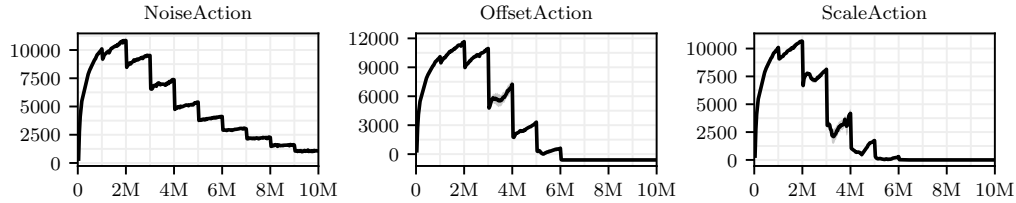


(b) Train returns.

Figure 20: Returns for the **sequential** training setup. A **SAC** agent is trained on the **HalfCheetah-v4** environment. **Continual adaptation:** The agent trains for  $T/10$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{10}\}$ . For **NoiseAction** tasks are determined by an action wrapper value in  $\{0.1, 0.2, \dots, 0.9\}$ , for **OffsetAction** in  $\{-0.1, -0.2, \dots, -0.9\}$ , for **ScaleAction** in  $\{0.9, 0.8, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. The modification is done on **action dimension 0**.



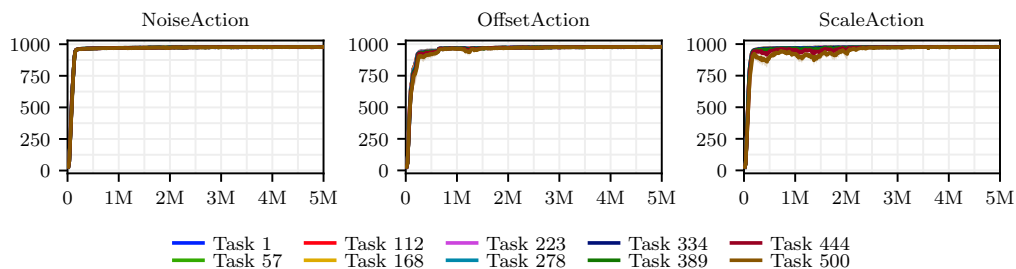
(a) Evaluation returns. Different tasks based on varied wrapper configuration values.



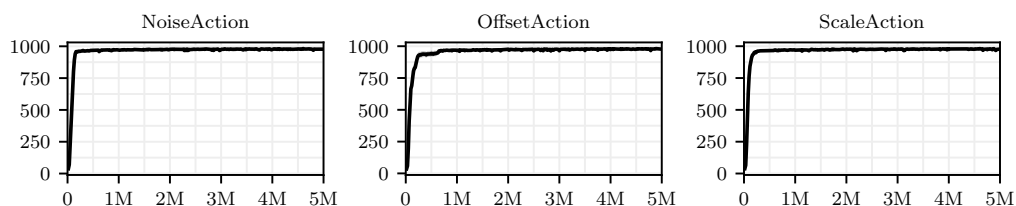
(b) Train returns.

Figure 21: Returns for the **sequential** training setup. A **SAC** agent is trained on the **HalfCheetah-v4** environment. **Continual adaptation:** The agent trains for  $T/10$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{10}\}$ . For **NoiseAction** tasks are determined by an action wrapper value in  $\{0.1, 0.2, \dots, 0.9\}$ , for **OffsetAction** in  $\{-0.1, -0.2, \dots, -0.9\}$ , for **ScaleAction** in  $\{0.9, 0.8, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. The modification is done on **all action dimensions**.



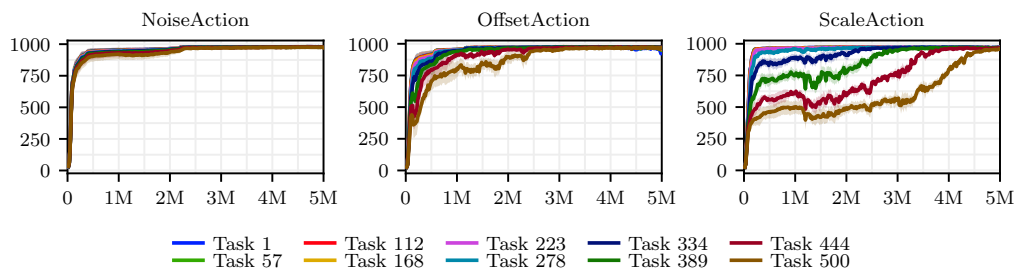


(a) Evaluation returns. Different tasks based on varied wrapper configuration values.

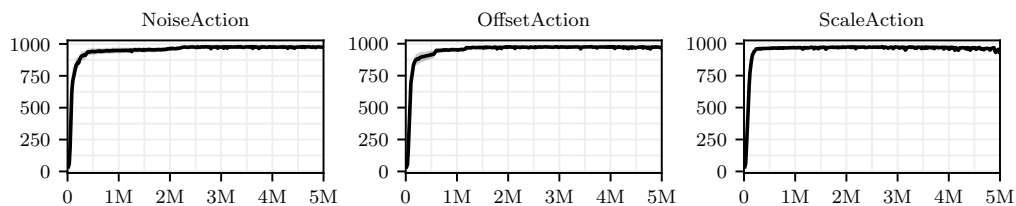


(b) Train returns.

Figure 22: Returns for the **sequential** training setup. A **SAC** agent is trained on the **walker-walk-v0** environment. **Continual slowly adaptation**: The agent trains for  $T/500$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then 499 times for  $T/500$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{500}\}$ . For **NoiseAction** tasks are determined by an action wrapper value in  $\{0.0, \dots, 0.9\}$ , for **OffsetAction** in  $\{0.0, \dots, -0.9\}$ , for **ScaleAction** in  $\{1.0, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. Only ten tasks are shown for better clarity. The modification is done on **action dimension 0**.

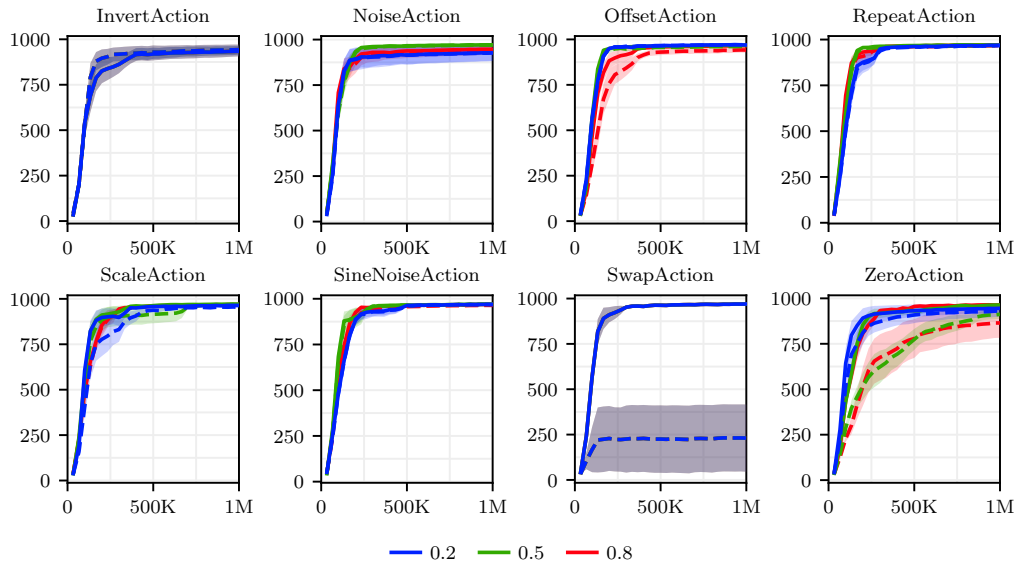


(a) Evaluation returns. Different tasks based on varied wrapper configuration values.

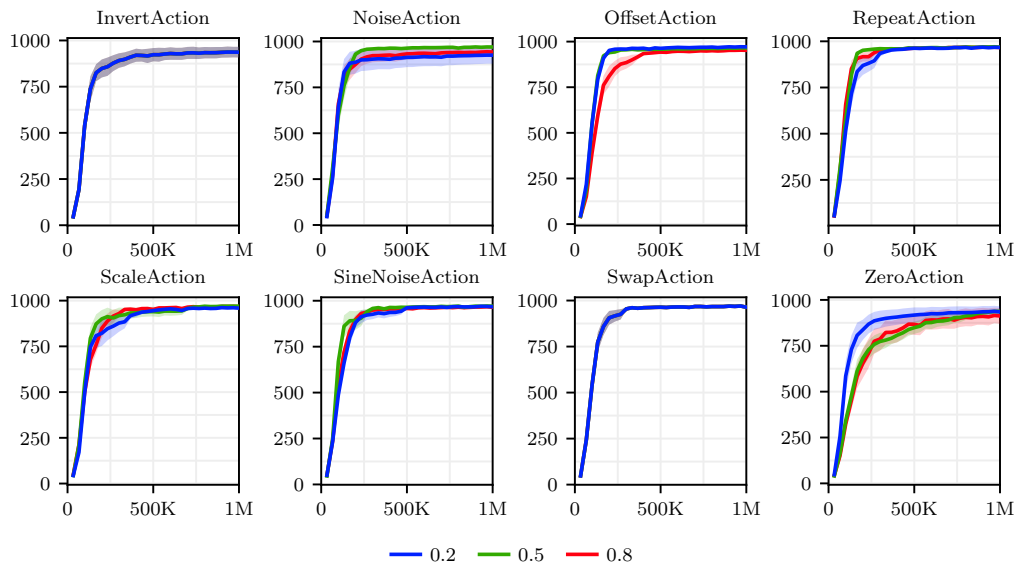


(b) Train returns.

Figure 23: Returns for the **sequential** training setup. A **SAC** agent is trained on the **walker-walk-v0** environment. **Continual slowly adaptation**: The agent trains for  $T/500$  steps on the unmodified action dynamic, task  $\mathcal{T}_1$ . Then 499 times for  $T/500$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value, tasks in  $\{\mathcal{T}_2, \mathcal{T}_3, \dots, \mathcal{T}_{500}\}$ . For **NoiseAction** tasks are determined by an action wrapper value in  $\{0.0, \dots, 0.9\}$ , for **OffsetAction** in  $\{0.0, \dots, -0.9\}$ , for **ScaleAction** in  $\{1.0, \dots, 0.1\}$ . The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the different modifications based on a varied wrapper value. Only ten tasks are shown for better clarity. The modification is done on **all action dimensions**.

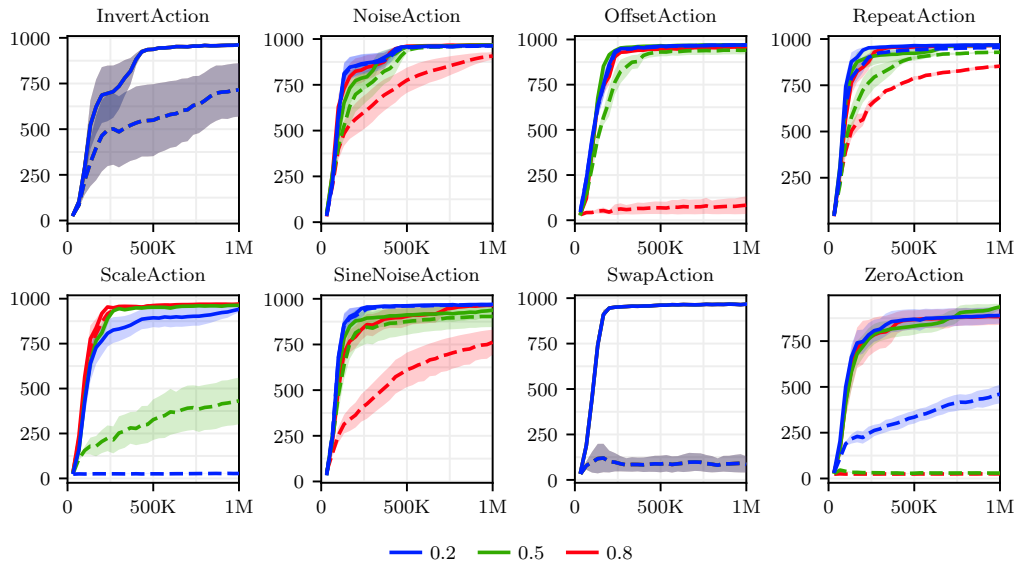


(a) Evaluation returns. Varied wrapper configuration value.

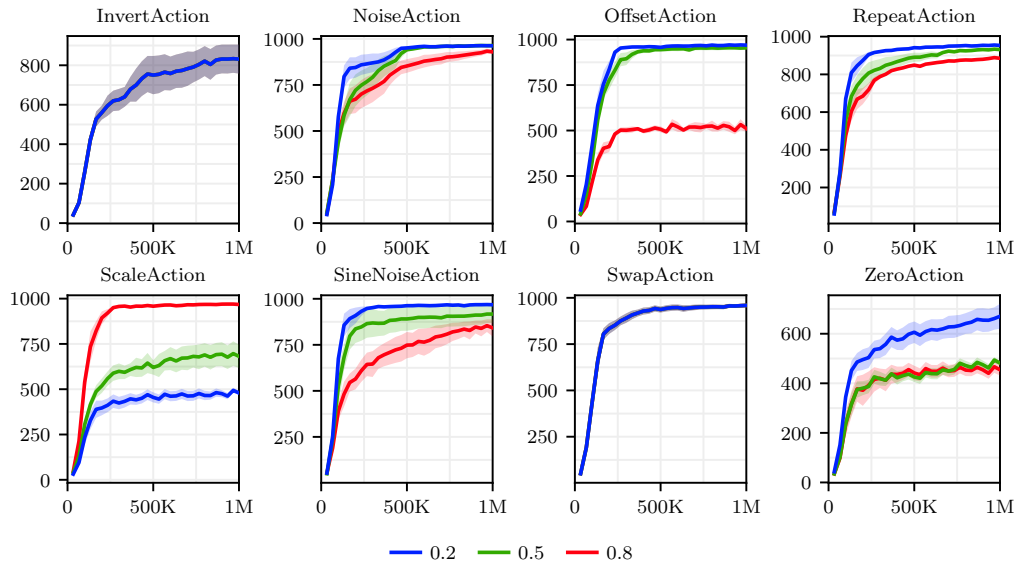


(b) Train returns. Varied wrapper configuration value.

Figure 24: Returns for the **multi-task** training setup. A SAC agent is trained on the `walker-walk-v0` environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.

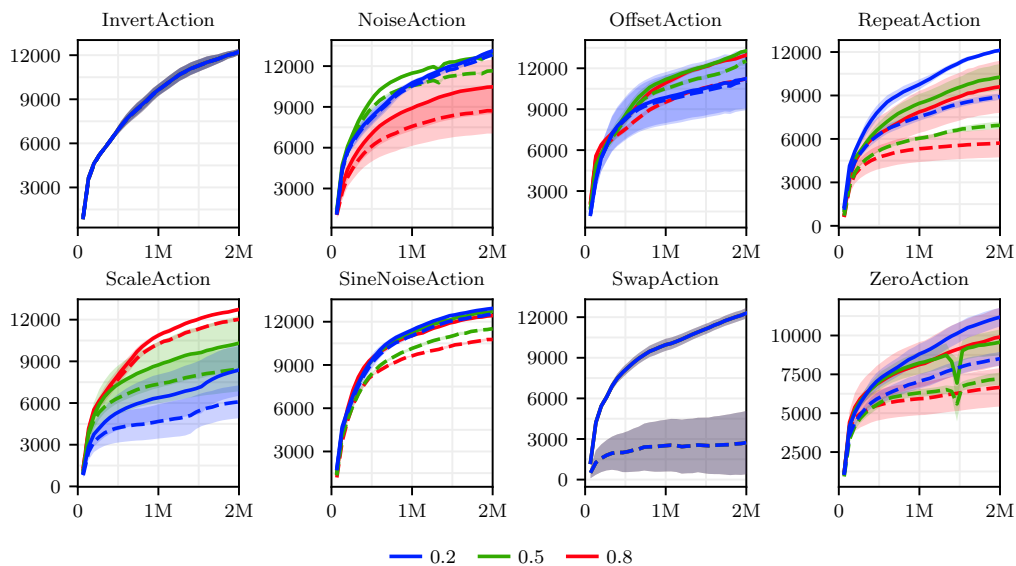


(a) Evaluation returns. Varied wrapper configuration value.

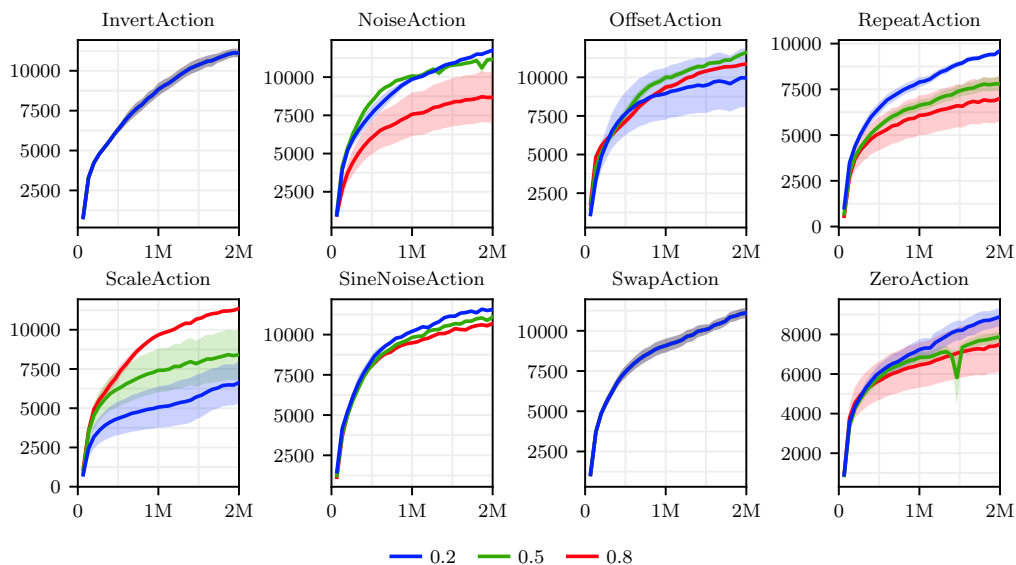


(b) Train returns. Varied wrapper configuration value.

Figure 25: Returns for the **multi-task** training setup. A SAC agent is trained on the `walker-walk-v0` environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

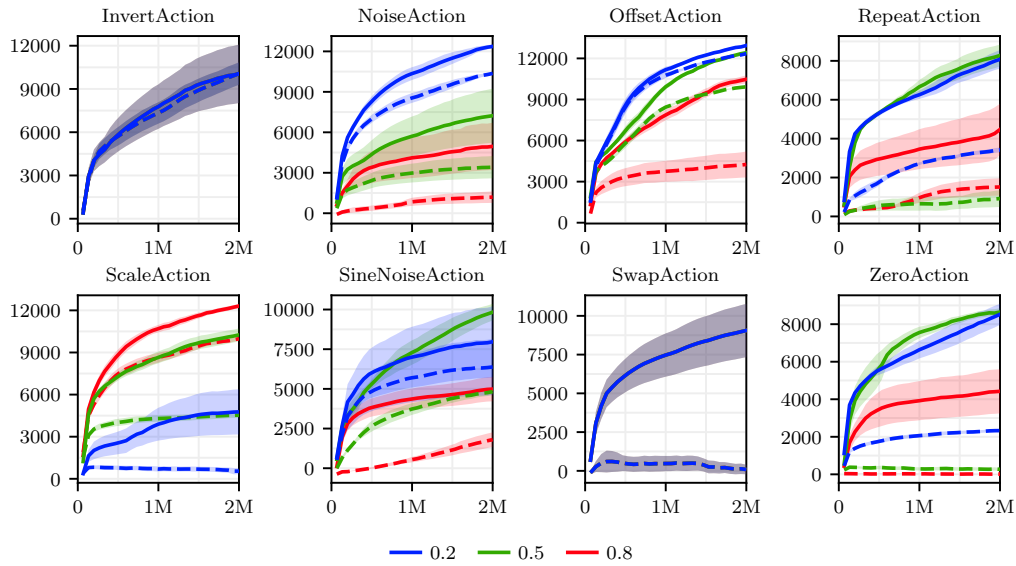


(a) Evaluation returns. Varied wrapper configuration value.

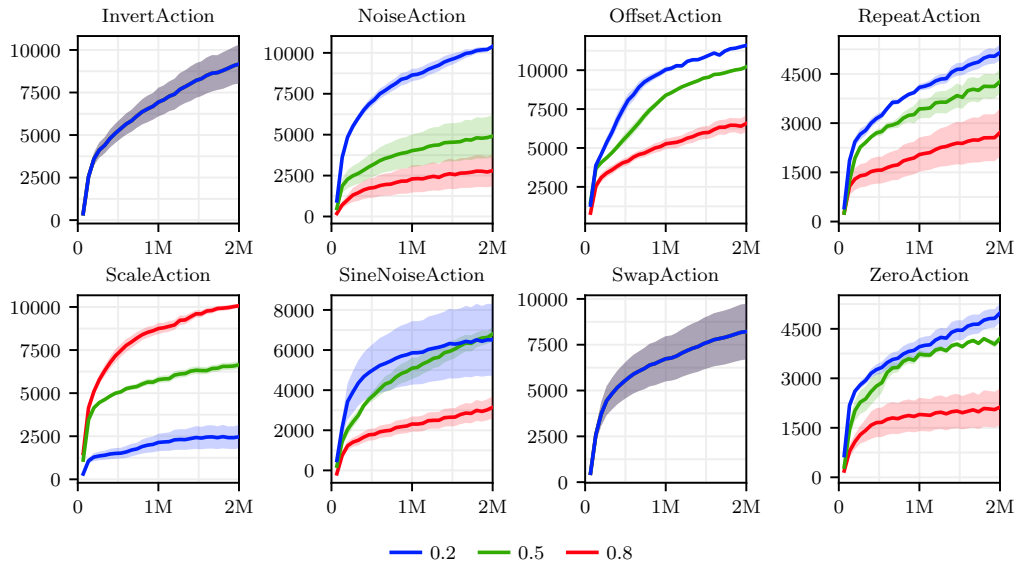


(b) Train returns. Varied wrapper configuration value.

Figure 26: Returns for the **multi-task** training setup. A SAC agent is trained on the HalfCheetah-v4 environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **action dimension 0**.



(a) Evaluation returns. Varied wrapper configuration value.



(b) Train returns. Varied wrapper configuration value.

Figure 27: Returns for the **multi-task** training setup. A SAC agent is trained on the HalfCheetah-v4 environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1. Colors represent the varied wrapper value. Solid lines represent the unmodified action dynamic, dashed lines the modified one. The modification is done on **all action dimensions**.

Table 4: Metrics for the **sequential** training setup. A **SAC** agent is trained on the `walker-walk-v0` environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Value	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
InvertAction	0	-	0.54 (0.01)	-0.01 (0.01)	0.0 (0.0)	0.92 (0.02)
	all	-	0.17 (0.06)	0.0 (0.0)	0.0 (0.0)	0.7 (0.12)
NoiseAction	0	0.2	1.0 (0.0)	0.96 (0.0)	0.01 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.96 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.92 (0.03)	0.01 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.96 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.85 (0.07)	0.01 (0.0)	0.0 (0.0)
		0.8	0.92 (0.02)	0.57 (0.04)	0.0 (0.0)	0.01 (0.01)
OffsetAction	0	0.2	1.0 (0.0)	0.96 (0.0)	0.01 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.95 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.8	0.99 (0.0)	0.45 (0.08)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.96 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.5	0.98 (0.01)	0.48 (0.05)	0.0 (0.0)	0.0 (0.0)
		0.8	0.45 (0.04)	0.05 (0.02)	0.0 (0.0)	0.32 (0.07)
RepeatAction	0	0.2	1.0 (0.0)	0.94 (0.02)	0.01 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.91 (0.04)	0.01 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.92 (0.03)	0.01 (0.0)	0.0 (0.0)
	all	0.2	0.99 (0.0)	0.8 (0.06)	0.01 (0.0)	0.0 (0.0)
		0.5	0.95 (0.02)	0.67 (0.04)	0.0 (0.0)	0.0 (0.0)
		0.8	0.74 (0.08)	0.58 (0.04)	0.0 (0.0)	0.12 (0.06)
ScaleAction	0	0.2	1.0 (0.0)	0.51 (0.05)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.9 (0.05)	0.01 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.96 (0.0)	0.01 (0.0)	0.0 (0.0)
	all	0.2	0.37 (0.05)	0.01 (0.0)	0.0 (0.0)	0.29 (0.1)
		0.5	0.72 (0.06)	0.26 (0.05)	0.0 (0.0)	0.08 (0.04)
		0.8	1.0 (0.0)	0.88 (0.08)	0.01 (0.0)	0.0 (0.0)
SineNoiseAction	0	0.2	1.0 (0.0)	0.96 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.95 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.87 (0.05)	0.01 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.95 (0.01)	0.01 (0.0)	0.0 (0.0)
		0.5	0.99 (0.0)	0.77 (0.07)	0.01 (0.0)	0.0 (0.0)
		0.8	0.81 (0.03)	0.44 (0.04)	0.0 (0.0)	0.0 (0.0)
SwapAction	0	-	0.88 (0.04)	0.06 (0.04)	0.0 (0.0)	0.24 (0.08)
	all	-	0.77 (0.06)	0.0 (0.0)	0.0 (0.0)	0.46 (0.12)
ZeroAction	0	0.2	1.0 (0.0)	0.78 (0.06)	0.01 (0.0)	0.0 (0.0)
		0.5	0.99 (0.0)	0.52 (0.05)	0.01 (0.0)	0.0 (0.0)
		0.8	0.96 (0.03)	0.34 (0.05)	0.0 (0.0)	0.01 (0.01)
	all	0.2	0.71 (0.01)	0.36 (0.03)	0.0 (0.0)	0.0 (0.0)
		0.5	0.5 (0.03)	0.07 (0.01)	0.0 (0.0)	0.06 (0.04)
		0.8	0.44 (0.02)	0.01 (0.0)	0.0 (0.0)	0.14 (0.03)

Table 5: Metrics for the **sequential** training setup. A **DrQ** agent is trained on the `walker-walk-v0(pixel)` environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Value	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
InvertAction	0	-	0.52 (0.01)	-0.0 (0.0)	0.0 (0.0)	0.73 (0.18)
	all	-	0.51 (0.01)	-0.0 (0.0)	0.0 (0.0)	0.92 (0.02)
NoiseAction	0	0.2	1.0 (0.01)	0.9 (0.02)	0.09 (0.04)	0.01 (0.01)
		0.5	0.99 (0.01)	0.78 (0.07)	0.13 (0.07)	0.0 (0.0)
		0.8	0.98 (0.01)	0.67 (0.13)	0.24 (0.13)	0.01 (0.01)
	all	0.2	0.98 (0.02)	0.9 (0.02)	0.06 (0.02)	0.01 (0.01)
		0.5	0.94 (0.02)	0.75 (0.03)	0.04 (0.02)	0.0 (0.0)
		0.8	0.88 (0.08)	0.52 (0.06)	0.06 (0.06)	0.05 (0.02)
OffsetAction	0	0.2	0.99 (0.02)	0.89 (0.02)	0.05 (0.03)	0.02 (0.01)
		0.5	0.98 (0.03)	0.67 (0.03)	0.04 (0.02)	0.05 (0.03)
		0.8	0.95 (0.02)	0.38 (0.02)	0.0 (0.0)	0.08 (0.03)
	all	0.2	0.99 (0.01)	0.8 (0.05)	0.04 (0.02)	0.0 (0.0)
		0.5	0.91 (0.03)	0.24 (0.05)	0.06 (0.04)	0.11 (0.06)
		0.8	0.29 (0.06)	0.03 (0.01)	0.0 (0.0)	0.88 (0.05)
RepeatAction	0	0.2	0.99 (0.01)	0.87 (0.04)	0.08 (0.02)	0.0 (0.0)
		0.5	0.98 (0.01)	0.89 (0.03)	0.06 (0.02)	0.0 (0.0)
		0.8	1.0 (0.02)	0.86 (0.05)	0.05 (0.03)	0.01 (0.01)
	all	0.2	0.94 (0.02)	0.77 (0.02)	0.02 (0.01)	0.0 (0.0)
		0.5	0.86 (0.02)	0.71 (0.03)	0.0 (0.0)	0.07 (0.02)
		0.8	0.85 (0.03)	0.67 (0.02)	0.0 (0.0)	0.09 (0.03)
ScaleAction	0	0.2	0.94 (0.03)	0.43 (0.06)	0.02 (0.01)	0.08 (0.04)
		0.5	0.99 (0.01)	0.77 (0.02)	0.03 (0.02)	0.03 (0.02)
		0.8	0.99 (0.03)	0.87 (0.02)	0.06 (0.02)	0.03 (0.03)
	all	0.2	0.04 (0.0)	0.0 (0.0)	0.0 (0.0)	0.96 (0.02)
		0.5	0.77 (0.03)	0.2 (0.03)	0.0 (0.0)	0.2 (0.04)
		0.8	0.98 (0.02)	0.74 (0.04)	0.07 (0.04)	0.02 (0.02)
SineNoiseAction	0	0.2	0.99 (0.02)	0.86 (0.03)	0.07 (0.02)	0.0 (0.0)
		0.5	1.0 (0.17)	0.82 (0.21)	0.08 (0.04)	0.0 (0.0)
		0.8	0.98 (0.01)	0.72 (0.03)	0.03 (0.02)	0.01 (0.01)
	all	0.2	1.0 (0.03)	0.8 (0.07)	0.14 (0.07)	0.0 (0.0)
		0.5	0.93 (0.01)	0.62 (0.03)	0.03 (0.02)	0.0 (0.0)
		0.8	0.81 (0.02)	0.41 (0.03)	0.02 (0.02)	0.06 (0.02)
SwapAction	0	-	0.26 (0.13)	0.14 (0.02)	0.0 (0.0)	0.74 (0.08)
	all	-	0.1 (0.03)	0.03 (0.02)	0.0 (0.0)	0.89 (0.05)
ZeroAction	0	0.2	0.97 (0.02)	0.62 (0.04)	0.08 (0.03)	0.0 (0.0)
		0.5	0.96 (0.13)	0.4 (0.1)	0.14 (0.09)	0.05 (0.03)
		0.8	0.96 (0.02)	0.27 (0.06)	0.01 (0.0)	0.04 (0.02)
	all	0.2	0.63 (0.02)	0.3 (0.04)	0.0 (0.0)	0.18 (0.04)
		0.5	0.07 (0.01)	0.06 (0.01)	0.0 (0.0)	0.92 (0.03)
		0.8	0.04 (0.0)	0.0 (0.0)	0.0 (0.0)	0.96 (0.03)



Table 6: Metrics for the **sequential** training setup. A **SAC** agent is trained on the `HalfCheetah-v4` environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Value	Avg. perf.	Fwd. trans.	Bwd. transfer	Forgetting
InvertAction	0	-	0.56 (0.01)	0.09 (0.01)	0.0 (0.0)	0.79 (0.02)
	all	-	0.4 (0.02)	0.01 (0.02)	0.0 (0.0)	1.02 (0.01)
NoiseAction	0	0.2	0.99 (0.02)	0.84 (0.02)	0.16 (0.02)	0.0 (0.0)
		0.5	0.98 (0.01)	0.83 (0.03)	0.09 (0.01)	0.0 (0.0)
		0.8	0.95 (0.01)	0.75 (0.02)	0.06 (0.02)	0.01 (0.01)
	all	0.2	0.95 (0.01)	0.72 (0.04)	0.03 (0.02)	0.0 (0.0)
		0.5	0.61 (0.01)	0.34 (0.03)	0.0 (0.0)	0.29 (0.03)
		0.8	0.39 (0.01)	0.2 (0.01)	0.0 (0.0)	0.51 (0.02)
OffsetAction	0	0.2	1.0 (0.03)	0.85 (0.02)	0.14 (0.02)	0.0 (0.0)
		0.5	1.0 (0.03)	0.81 (0.02)	0.09 (0.03)	0.0 (0.0)
		0.8	0.98 (0.01)	0.69 (0.04)	0.01 (0.01)	0.02 (0.01)
	all	0.2	0.94 (0.02)	0.62 (0.09)	0.0 (0.0)	0.04 (0.01)
		0.5	0.82 (0.04)	0.26 (0.04)	0.0 (0.0)	0.22 (0.06)
		0.8	0.01 (0.01)	0.09 (0.02)	0.0 (0.0)	0.97 (0.02)
RepeatAction	0	0.2	0.9 (0.02)	0.56 (0.03)	0.0 (0.0)	0.02 (0.0)
		0.5	0.8 (0.02)	0.39 (0.02)	0.0 (0.0)	0.12 (0.03)
		0.8	0.77 (0.02)	0.33 (0.02)	0.0 (0.0)	0.16 (0.02)
	all	0.2	0.51 (0.02)	0.25 (0.03)	0.0 (0.0)	0.39 (0.03)
		0.5	0.32 (0.01)	0.12 (0.03)	0.0 (0.0)	0.63 (0.02)
		0.8	0.25 (0.02)	0.08 (0.02)	0.0 (0.0)	0.75 (0.08)
ScaleAction	0	0.2	0.47 (0.06)	0.4 (0.01)	0.0 (0.0)	0.66 (0.11)
		0.5	0.94 (0.02)	0.74 (0.02)	0.02 (0.01)	0.04 (0.02)
		0.8	0.98 (0.02)	0.83 (0.02)	0.12 (0.02)	0.0 (0.0)
	all	0.2	0.0 (0.0)	0.01 (0.01)	0.0 (0.0)	1.03 (0.02)
		0.5	0.01 (0.03)	0.21 (0.01)	0.0 (0.0)	1.01 (0.02)
		0.8	0.95 (0.01)	0.76 (0.02)	0.01 (0.01)	0.02 (0.01)
SineNoiseAction	0	0.2	0.99 (0.02)	0.84 (0.02)	0.15 (0.02)	0.0 (0.0)
		0.5	0.95 (0.02)	0.8 (0.02)	0.12 (0.03)	0.0 (0.0)
		0.8	0.94 (0.02)	0.76 (0.04)	0.06 (0.02)	0.0 (0.0)
	all	0.2	0.94 (0.01)	0.65 (0.06)	0.01 (0.0)	0.0 (0.0)
		0.5	0.62 (0.03)	0.24 (0.02)	0.0 (0.0)	0.26 (0.07)
		0.8	0.37 (0.01)	0.08 (0.01)	0.0 (0.0)	0.52 (0.02)
SwapAction	0	-	0.57 (0.03)	0.07 (0.03)	0.0 (0.0)	0.81 (0.08)
	all	-	0.5 (0.02)	0.02 (0.02)	0.0 (0.0)	0.98 (0.03)
ZeroAction	0	0.2	0.88 (0.01)	0.59 (0.02)	0.01 (0.01)	0.03 (0.02)
		0.5	0.72 (0.03)	0.39 (0.01)	0.0 (0.0)	0.24 (0.06)
		0.8	0.55 (0.06)	0.32 (0.01)	0.0 (0.0)	0.51 (0.13)
	all	0.2	0.33 (0.02)	0.21 (0.01)	0.0 (0.0)	0.58 (0.03)
		0.5	0.05 (0.01)	0.03 (0.0)	0.0 (0.0)	0.94 (0.03)
		0.8	0.0 (0.01)	0.0 (0.0)	0.0 (0.0)	1.0 (0.02)

Table 7: Metrics for the **sequential** training setup. A **SAC** agent is trained with **HER** on the goal-conditioned **FetchReach-v2** environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Value	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
InvertAction	0	-	0.38 (0.07)	-0.04 (0.02)	0.0 (0.0)	0.88 (0.1)
	all	-	0.21 (0.08)	-0.02 (0.02)	0.0 (0.0)	1.0 (0.0)
NoiseAction	0	0.2	1.0 (0.0)	0.94 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.98 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.8	0.98 (0.01)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	1.0 (0.0)	0.0 (0.0)	0.0 (0.0)
		0.5	0.99 (0.01)	0.92 (0.04)	0.0 (0.0)	0.0 (0.0)
		0.8	0.94 (0.03)	0.84 (0.04)	0.0 (0.0)	0.0 (0.0)
OffsetAction	0	0.2	1.0 (0.0)	1.0 (0.0)	0.0 (0.0)	0.0 (0.0)
		0.5	0.99 (0.01)	1.0 (0.0)	0.0 (0.0)	0.0 (0.0)
		0.8	0.94 (0.04)	0.98 (0.02)	0.0 (0.0)	0.02 (0.02)
	all	0.2	0.99 (0.01)	1.0 (0.0)	0.0 (0.0)	0.02 (0.02)
		0.5	0.94 (0.05)	1.0 (0.0)	0.0 (0.0)	0.12 (0.1)
		0.8	0.97 (0.02)	0.76 (0.08)	0.0 (0.0)	0.0 (0.0)
RepeatAction	0	0.2	1.0 (0.0)	0.94 (0.04)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.92 (0.06)	0.0 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.8 (0.06)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.92 (0.06)	0.0 (0.0)	0.0 (0.0)
		0.5	0.99 (0.01)	0.62 (0.12)	0.0 (0.0)	0.0 (0.0)
		0.8	0.95 (0.04)	0.52 (0.02)	0.0 (0.0)	0.02 (0.02)
ScaleAction	0	0.2	0.98 (0.02)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.98 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
SineNoiseAction	0	0.2	1.0 (0.0)	0.94 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.5	0.99 (0.01)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.8	0.99 (0.01)	0.78 (0.07)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	1.0 (0.0)	0.0 (0.0)	0.0 (0.0)
		0.5	0.98 (0.01)	0.9 (0.04)	0.0 (0.0)	0.0 (0.0)
		0.8	0.87 (0.01)	0.68 (0.11)	0.0 (0.0)	0.0 (0.0)
SwapAction	0	-	0.61 (0.17)	0.0 (0.04)	0.0 (0.0)	0.54 (0.22)
	all	-	0.87 (0.13)	0.0 (0.03)	0.0 (0.0)	0.2 (0.2)
ZeroAction	0	0.2	0.98 (0.02)	0.96 (0.02)	0.0 (0.0)	0.02 (0.02)
		0.5	1.0 (0.0)	0.94 (0.04)	0.0 (0.0)	0.0 (0.0)
		0.8	0.88 (0.03)	0.7 (0.08)	0.0 (0.0)	0.0 (0.0)
	all	0.2	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.5	1.0 (0.0)	0.96 (0.02)	0.0 (0.0)	0.0 (0.0)
		0.8	0.82 (0.02)	0.58 (0.02)	0.0 (0.0)	0.04 (0.02)

Table 8: Metrics for the **sequential** training setup. A **SAC** agent is trained with **HER** on the goal-conditioned **FetchPush-v2** environment. **One dynamics switch**: The agent trains for  $T/2$  steps on the unmodified action dynamic, and then for  $T/2$  steps on the modified action dynamic. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Value	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
InvertAction	0	-	0.24 (0.12)	0.1 (0.05)	0.0 (0.0)	0.59 (0.2)
	all	-	0.26 (0.1)	0.08 (0.04)	0.0 (0.0)	0.55 (0.19)
NoiseAction	0	0.2	1.0 (0.24)	0.67 (0.25)	0.36 (0.23)	0.0 (0.0)
		0.5	0.98 (0.3)	0.88 (0.34)	0.25 (0.15)	0.03 (0.03)
		0.8	0.96 (0.33)	0.79 (0.32)	0.21 (0.21)	0.11 (0.11)
	all	0.2	1.0 (0.01)	0.78 (0.15)	0.2 (0.16)	0.0 (0.0)
		0.5	0.98 (0.24)	0.85 (0.18)	0.25 (0.19)	0.28 (0.24)
		0.8	0.92 (0.16)	0.71 (0.18)	0.22 (0.16)	0.17 (0.17)
OffsetAction	0	0.2	0.65 (0.22)	0.92 (0.05)	0.0 (0.0)	0.33 (0.2)
		0.5	0.78 (0.14)	0.84 (0.04)	0.0 (0.0)	0.29 (0.14)
		0.8	0.6 (0.16)	0.31 (0.06)	0.0 (0.0)	0.33 (0.18)
	all	0.2	0.43 (0.21)	0.92 (0.06)	0.0 (0.0)	0.53 (0.22)
		0.5	0.54 (0.2)	0.74 (0.06)	0.0 (0.0)	0.47 (0.18)
		0.8	0.41 (0.16)	0.18 (0.07)	0.0 (0.0)	0.47 (0.21)
RepeatAction	0	0.2	0.98 (0.31)	0.66 (0.31)	0.38 (0.26)	0.09 (0.06)
		0.5	0.54 (0.21)	0.95 (0.1)	0.02 (0.02)	0.54 (0.24)
		0.8	0.99 (0.22)	0.88 (0.12)	0.12 (0.08)	0.1 (0.07)
	all	0.2	0.98 (0.3)	0.67 (0.35)	0.33 (0.26)	0.07 (0.04)
		0.5	0.77 (0.2)	0.88 (0.05)	0.02 (0.02)	0.24 (0.17)
		0.8	0.94 (0.16)	0.64 (0.09)	0.08 (0.05)	0.08 (0.05)
ScaleAction	0	0.2	0.57 (0.16)	0.31 (0.06)	0.0 (0.0)	0.31 (0.19)
		0.5	0.62 (0.16)	0.9 (0.04)	0.0 (0.0)	0.37 (0.17)
		0.8	0.82 (0.14)	0.94 (0.04)	0.0 (0.0)	0.22 (0.18)
	all	0.2	0.44 (0.14)	0.35 (0.08)	0.0 (0.0)	0.33 (0.2)
		0.5	0.44 (0.23)	0.82 (0.05)	0.0 (0.0)	0.55 (0.22)
		0.8	0.46 (0.18)	0.92 (0.03)	0.0 (0.0)	0.55 (0.17)
SineNoiseAction	0	0.2	0.89 (0.34)	0.96 (0.36)	0.07 (0.07)	0.11 (0.07)
		0.5	0.88 (0.41)	0.96 (0.41)	0.08 (0.08)	0.23 (0.14)
		0.8	0.98 (0.31)	0.72 (0.3)	0.25 (0.18)	0.06 (0.06)
	all	0.2	0.98 (0.08)	0.83 (0.13)	0.19 (0.19)	0.06 (0.06)
		0.5	0.95 (0.18)	0.79 (0.22)	0.26 (0.21)	0.19 (0.19)
		0.8	0.92 (0.17)	0.44 (0.15)	0.07 (0.05)	0.02 (0.02)
SwapAction	0	-	0.51 (0.16)	0.2 (0.08)	0.0 (0.0)	0.43 (0.15)
	all	-	0.67 (0.1)	0.24 (0.12)	0.0 (0.0)	0.24 (0.06)
ZeroAction	0	0.2	0.98 (0.31)	0.61 (0.32)	0.33 (0.23)	0.03 (0.03)
		0.5	0.85 (0.16)	0.58 (0.1)	0.1 (0.07)	0.12 (0.08)
		0.8	0.62 (0.02)	0.17 (0.04)	0.17 (0.1)	0.02 (0.02)
	all	0.2	0.98 (0.45)	0.87 (0.46)	0.44 (0.34)	0.44 (0.38)
		0.5	0.68 (0.14)	0.49 (0.13)	0.1 (0.07)	0.22 (0.16)
		0.8	0.56 (0.03)	0.13 (0.04)	0.19 (0.1)	0.02 (0.02)

Table 9: Metrics for the **sequential** training setup. A **SAC** agent is trained on the `walker-walk-v0` environment. **Continual adaptation**: The agent trains for  $T/10$  steps on the unmodified action dynamic, and then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value (increasing for `NoiseAction` and `OffsetAction`, decreasing for `ScaleAction`). The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
NoiseAction	0	1.0 (0.0)	0.19 (0.0)	0.0 (0.0)	0.0 (0.0)
	all	0.97 (0.0)	0.19 (0.0)	0.0 (0.0)	0.0 (0.0)
OffsetAction	0	0.71 (0.05)	0.19 (0.0)	0.04 (0.02)	0.09 (0.01)
	all	0.04 (0.0)	0.11 (0.0)	0.0 (0.0)	0.11 (0.0)
ScaleAction	0	0.98 (0.0)	0.19 (0.0)	0.0 (0.0)	0.01 (0.0)
	all	0.03 (0.0)	0.1 (0.0)	0.0 (0.0)	0.11 (0.0)

Table 10: Metrics for the **sequential** training setup. A **SAC** agent is trained on the `HalfCheetah-v4` environment. **Continual adaptation**: The agent trains for  $T/10$  steps on the unmodified action dynamic, and then nine times for  $T/10$  steps on modified action dynamics with increasing difficulty, determined by the wrapper value (increasing for `NoiseAction` and `OffsetAction`, decreasing for `ScaleAction`). The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

Wrapper	Dim.	Avg. perf.	Fwd. trans.	Bwd. trans.	Forgetting
NoiseAction	0	0.83 (0.01)	0.18 (0.0)	0.01 (0.0)	0.02 (0.0)
	all	0.34 (0.0)	0.12 (0.0)	0.0 (0.0)	0.06 (0.0)
OffsetAction	0	0.36 (0.02)	0.15 (0.0)	0.04 (0.01)	0.11 (0.01)
	all	-0.03 (0.0)	0.08 (0.0)	0.01 (0.0)	0.08 (0.0)
ScaleAction	0	0.81 (0.03)	0.17 (0.0)	0.04 (0.0)	0.06 (0.01)
	all	-0.02 (0.0)	0.06 (0.0)	0.01 (0.0)	0.08 (0.0)

Table 11: Metrics for the **multi-task** training setup. A **SAC** agent is trained on the `walker-walk-v0` environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

<b>Wrapper</b>	<b>Dim.</b>	<b>Value</b>	<b>Avg. perf.</b>	<b>Parallel trans.</b>
InvertAction	0	-	0.99 (0.03)	-0.02 (0.03)
	all	-	0.86 (0.08)	-0.14 (0.08)
NoiseAction	0	0.2	1.0 (0.05)	-0.04 (0.05)
		0.5	1.0 (0.0)	0.01 (0.01)
		0.8	1.0 (0.02)	-0.0 (0.0)
	all	0.2	1.0 (0.01)	-0.0 (0.01)
		0.5	1.0 (0.0)	-0.0 (0.0)
		0.8	0.97 (0.01)	-0.02 (0.01)
OffsetAction	0	0.2	0.98 (0.01)	-0.0 (0.02)
		0.5	1.0 (0.0)	0.02 (0.01)
		0.8	0.99 (0.01)	-0.0 (0.01)
	all	0.2	1.0 (0.0)	-0.0 (0.0)
		0.5	0.98 (0.01)	-0.0 (0.01)
		0.8	0.55 (0.03)	-0.28 (0.08)
RepeatAction	0	0.2	1.0 (0.0)	-0.0 (0.0)
		0.5	0.99 (0.02)	-0.02 (0.02)
		0.8	1.0 (0.0)	0.0 (0.0)
	all	0.2	0.99 (0.0)	-0.0 (0.0)
		0.5	0.98 (0.0)	-0.0 (0.01)
		0.8	0.95 (0.0)	-0.03 (0.0)
ScaleAction	0	0.2	1.0 (0.0)	-0.0 (0.0)
		0.5	1.0 (0.0)	0.01 (0.01)
		0.8	1.0 (0.0)	-0.0 (0.0)
	all	0.2	0.51 (0.01)	-0.08 (0.02)
		0.5	0.71 (0.08)	-0.17 (0.07)
		0.8	1.0 (0.0)	0.01 (0.01)
SineNoiseAction	0	0.2	1.0 (0.0)	-0.0 (0.0)
		0.5	1.0 (0.0)	0.0 (0.0)
		0.8	1.0 (0.0)	-0.0 (0.0)
	all	0.2	1.0 (0.0)	0.0 (0.0)
		0.5	0.98 (0.05)	-0.05 (0.05)
		0.8	0.9 (0.04)	-0.08 (0.04)
SwapAction	0	-	0.62 (0.1)	-0.0 (0.01)
	all	-	0.55 (0.02)	0.02 (0.02)
ZeroAction	0	0.2	0.99 (0.03)	-0.03 (0.03)
		0.5	0.97 (0.01)	0.0 (0.03)
		0.8	0.95 (0.04)	-0.01 (0.06)
	all	0.2	0.76 (0.06)	-0.16 (0.06)
		0.5	0.52 (0.01)	-0.08 (0.01)
		0.8	0.52 (0.02)	-0.04 (0.02)

Table 12: Metrics for the **multi-task** training setup. A **SAC** agent is trained on the **HalfCheetah-v4** environment. The agent trains on two tasks resulting from the unmodified action dynamic and one modification of such. The modification is configured with the action wrapper, the action dimension and if applicable a wrapper specific value, see Section 4 and Table 1.

<b>Wrapper</b>	<b>Dim.</b>	<b>Value</b>	<b>Avg. perf.</b>	<b>Parallel trans.</b>
InvertAction	0	-	1.0 (0.02)	-0.06 (0.01)
	all	-	0.99 (0.12)	-0.16 (0.02)
NoiseAction	0	0.2	0.98 (0.02)	0.0 (0.03)
		0.5	0.95 (0.01)	-0.02 (0.01)
		0.8	0.92 (0.18)	-0.11 (0.18)
	all	0.2	0.92 (0.01)	0.02 (0.07)
		0.5	0.73 (0.19)	-0.34 (0.11)
		0.8	0.62 (0.22)	-0.42 (0.08)
OffsetAction	0	0.2	1.0 (0.19)	-0.13 (0.16)
		0.5	0.97 (0.02)	0.01 (0.04)
		0.8	0.93 (0.01)	-0.02 (0.02)
	all	0.2	0.98 (0.01)	-0.01 (0.0)
		0.5	0.9 (0.01)	-0.04 (0.02)
		0.8	0.7 (0.06)	-0.2 (0.05)
RepeatAction	0	0.2	0.87 (0.01)	-0.02 (0.06)
		0.5	0.84 (0.04)	-0.09 (0.1)
		0.8	0.8 (0.14)	-0.21 (0.1)
	all	0.2	0.71 (0.03)	-0.18 (0.03)
		0.5	0.56 (0.05)	-0.22 (0.03)
		0.8	0.67 (0.18)	-0.38 (0.07)
ScaleAction	0	0.2	0.86 (0.18)	-0.26 (0.08)
		0.5	0.91 (0.18)	-0.17 (0.13)
		0.8	0.97 (0.01)	-0.06 (0.04)
	all	0.2	0.56 (0.14)	-0.34 (0.06)
		0.5	0.72 (0.02)	-0.11 (0.02)
		0.8	0.91 (0.01)	-0.05 (0.02)
SineNoiseAction	0	0.2	0.98 (0.01)	-0.02 (0.02)
		0.5	0.96 (0.0)	0.01 (0.02)
		0.8	0.93 (0.01)	-0.03 (0.02)
	all	0.2	0.93 (0.26)	-0.36 (0.15)
		0.5	0.74 (0.03)	-0.18 (0.02)
		0.8	0.67 (0.12)	-0.39 (0.05)
SwapAction	0	-	0.61 (0.09)	-0.04 (0.02)
	all	-	0.51 (0.11)	-0.16 (0.07)
ZeroAction	0	0.2	0.88 (0.05)	-0.07 (0.06)
		0.5	0.88 (0.04)	-0.21 (0.02)
		0.8	0.85 (0.16)	-0.21 (0.1)
	all	0.2	0.64 (0.03)	-0.19 (0.02)
		0.5	0.51 (0.01)	-0.2 (0.01)
		0.8	0.5 (0.13)	-0.34 (0.04)