# MaDi: Learning to Mask Distractions for Generalization in Visual Deep Reinforcement Learning

Bram Grooten[1], Tristan Tomilin[1], Gautham Vasan[2], Matthew E. Taylor[2,3],
A. Rupam Mahmood[2,3], Meng Fang[4], Mykola Pechenizkiy[1], Decebal Constantin Mocanu[5,1]

[1]Eindhoven University of Technology  [2]University of Alberta  [3]Alberta Machine Intelligence Institute (Amii)
[4]University of Liverpool  [5]University of Luxembourg

## ABSTRACT

The visual world provides an abundance of information, but many input pixels received by agents often contain distracting stimuli. Autonomous agents need the ability to distinguish useful information from task-irrelevant perceptions, enabling them to generalize to unseen environments with new distractions. Existing works approach this problem using data augmentation or large auxiliary networks with additional loss functions. We introduce *MaDi*, a novel algorithm that learns to **ma**sk **di**stractions by the reward signal only. In MaDi, the conventional actor-critic structure of deep reinforcement learning agents is complemented by a small third sibling, the Masker. This lightweight neural network generates a mask to determine what the actor and critic will receive, such that they can focus on learning the task. The masks are created dynamically, depending on the current input. We run experiments on the DeepMind Control Generalization Benchmark, the Distracting Control Suite, and a real UR5 Robotic Arm. Our algorithm improves the agent's focus with useful masks, while its efficient Masker network only adds 0.2% more parameters to the original structure, in contrast to previous work. MaDi consistently achieves generalization results better than or competitive to state-of-the-art methods.[1]

## KEYWORDS

Deep Reinforcement Learning, Generalization, Robotics

**ACM Reference Format:**
Bram Grooten[1], Tristan Tomilin[1], Gautham Vasan[2], Matthew E. Taylor[2,3],, A. Rupam Mahmood[2,3], Meng Fang[4], Mykola Pechenizkiy[1], Decebal Constantin Mocanu[5,1] . 2024. MaDi: Learning to Mask Distractions for Generalization in Visual Deep Reinforcement Learning. In *Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024), Auckland, New Zealand, May 6 − 10, 2024*, IFAAMAS, 23 pages.

## 1 INTRODUCTION

Deep reinforcement learning (RL) has achieved remarkable success in a variety of complex tasks such as game playing [32, 37], robotics [2, 10, 17], nuclear fusion [7], and autonomous navigation [31, 60]. However, one of the major challenges faced by RL agents is their limited ability to generalize to unseen environments, particularly in the presence of distracting visual noise, such as a video playing in the background [22, 39]. These distractions can lead to significant degradation in the performance of deep RL agents, thereby hindering their applicability in the real world. To address this, we propose a novel algorithm, **Ma**sking **Di**stractions, which learns to filter out task-irrelevant visuals, enhancing generalization capabilities.

The key idea behind MaDi is to supplement the conventional actor-critic architecture with a third lightweight component, the Masker (see Figure 1). This small neural network generates a mask that dims the irrelevant pixels, allowing the actor and critic to focus on learning the task at hand without getting too distracted. Unlike previous approaches that have attempted to address this issue [4, 21, 22], our method increases generalization performance while introducing minimal overhead in terms of model parameters, thus preserving the efficiency of the original architecture.

Furthermore, no additional loss function is necessary for the Masker to optimize its parameters. To ensure that the Masker maintains visibility of the task-relevant pixels, it is trained on the critic's loss function. The Masker and critic networks are aligned in their objective, as pixels that are essential to determine the value of an observation should not be hidden. Figure 1 shows an example of a mask, visualizing the output produced by the Masker network corresponding to the current input frame. The Masker is able to learn such precise segmentations without any additional labels, bounding boxes, or other annotations. The reward alone is enough.

To evaluate the effectiveness of MaDi, we conduct experiments on multiple environments from three benchmarks: the DeepMind Control Generalization Benchmark [22], the Distracting Control Suite [39], and a real UR5 Robotic Arm for which we design a novel generalization experiment with visual distractions. Our results demonstrate that MaDi significantly improves the agent's ability to focus on relevant visual information by generating helpful masks, leading to enhanced generalization performance. Furthermore, MaDi achieves state-of-the-art performance on many environments, surpassing well-known methods in vision-based reinforcement learning [4, 18, 21, 22, 28, 50].

Our main contributions are:

- We introduce a novel algorithm, MaDi, which supplements the standard actor-critic architecture of deep RL agents with a lightweight Masker. This network learns to focus on the task-relevant pixels solely from the reward signal.
- We present a comprehensive set of experiments on the DeepMind Control Generalization Benchmark and the Distracting Control Suite. MaDi consistently achieves state-of-the-art or competitive generalization performance.
- We test MaDi on a physical robot, demonstrating that our algorithm increases the performance of the UR5 Robotic Arm in a challenging VisualReacher task, even when distracting videos are playing in the background.

---

[1]Code: github.com/bramgrooten/mask-distractions and video: youtu.be/2oImF0h1k48. Corresponding author: b.j.grooten@tue.nl
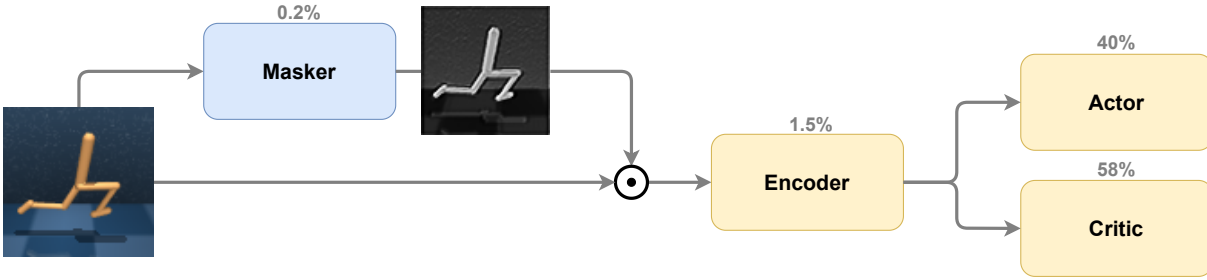
**Figure 1: In the MaDi architecture, the Masker produces a soft mask (values between 0 and 1) for each frame, which subsequently gets multiplied element-wise with the observation. The encoder is a shared ConvNet updated only by the critic loss, which is also the objective function for the Masker. We show rounded percentages for the number of parameters used in the `walker-walk` environment. The actor and critic contain most of it ($\sim 98\%$ together) as they consist of multiple fully-connected layers.**

The paper is structured as follows: Section 2 reviews related work, Section 3 formalizes our mathematical framework. Our algorithm MaDi is detailed in Section 4. We present simulation results in Section 5 and robotic experiments in Section 6. Section 7 concludes.

## 2 RELATED WORK

The problem of generalization in deep reinforcement learning has been an active area of research, with several approaches proposed to tackle the challenge of visual distractions. In this section, we review the most relevant literature, highlighting the differences between our proposed MaDi method and existing approaches.

*Generalization in RL.* In reinforcement learning, generalization refers to an agent's ability to perform well on unseen environments or tasks [27]. This can be challenging, as RL is prone to overfit to the training environment [6, 12, 20, 57]. Several works have focused on improving generalization capabilities by employing techniques such as domain adaptation [49], domain randomization [2, 42], meta-learning [9, 46], contrastive learning [1, 29], imitation learning [11], bisimulation metrics [13, 56], and data augmentation [22, 28, 34, 47, 50, 53]. Even using a ResNet [23] pretrained on ImageNet [36] as an encoder can improve generalization [54].

*Visual Learning in RL.* Learning tasks from visual input, i.e., image-based or vision-based deep RL, is typically more demanding than learning from direct features in a vector. DQN [32] was the first to learn Atari games at human-level performance directly from pixels. However, it has been shown that these algorithms can be quite brittle to changes in the environment, as altering a few pixel values can significantly decrease DQN's performance [33, 58]. Using data augmentation proved to be the key in visual RL. DrQ [50] and RAD [28] use light augmentations such as random shifts or crops of the observation to increase the algorithm's robustness.

*Distractions in RL.* Several approaches have been proposed to deal with the presence of task-irrelevant noise and distractions in reinforcement learning environments. Automatic Noise Filtering (ANF) [16] works on noisy environments that provide states as feature vectors, such as the MuJoCo Gym suite [5, 43]. We focus our work on RL agents that need to learn from image-based observations, like the pixel-wrapped DeepMind Control Suite [41]. Two benchmarks that we use are extensions of this suite.

Several works [4, 21, 22, 53, 54] have tried to tackle the DeepMind Control Generalization Benchmark [22] and the Distracting Control Suite [39]. Many of these methods use stronger[2] data augmentations than the light shifting and cropping of DrQ and RAD. Usually, they apply one of two favored augmentation techniques: a randomly initialized convolution layer (`conv` augmentation) or overlaying the observation with random images from a large dataset, such as Places365 [59] (`overlay` augmentation).

*Masking in Visual RL.* There exist a few works that aim to improve the generalization ability of RL agents by masking parts of the input. Yu et al. [51] randomly mask parts of the inputs and use an auxiliary loss to reconstruct these pixels. SGQN [4] and the recent InfoGating [44] apply more targeted masking similar to MaDi. InfoGating has only experimented on offline RL, and it uses a large U-Net [35] to determine the appropriate masks, while MaDi uses a much smaller 3-layer convolutional neural network.

*Baselines.* We select the following set of six baselines for our experiments, as these focus on online RL and do not use any pre-trained models:

- *Soft Actor-Critic* [18, **SAC**] is an off-policy actor-critic algorithm that optimizes the trade-off between exploration and exploitation by automatically tuning a temperature parameter for entropy regularization.
- *Data-regularized Q-learning* [50, **DrQ**] focuses on making Q-learning more stable and sample efficient by shifting the observations by a few pixels in a random direction.
- *RL with Augmented Data* [28, **RAD**] improves the data efficiency in visual RL by randomly cropping the images.
- *Soft Data Augmentation* [22, **SODA**] applies data augmentation in an auxiliary task that tries to minimize the distance of augmented and non-augmented images in its feature space.
- *Stabilized Value Estimation under Augmentation* [21, **SVEA**] stabilizes learning by using augmentation solely in the critic. It combines clean and augmented data in every batch used for a critic update. The actor only sees clean data.
- *Saliency Guided Q-Networks* [4, **SGQN**] is perhaps closest to our work, as it also uses masks to benefit learning. Its masks are not applied at the start of the architecture, but are learned

---

[2]By *stronger*, we mean augmentations that alter an image significantly more.

by a third component after the encoder. This auxiliary model minimizes the difference between its masks and other masks generated by a saliency metric. By computing the gradient of the Q-function with respect to the input pixels, this saliency metric determines which pixels are important for the agent. A hyperparameter `sgqn_quantile` (often set to $95\% - 98\%$) determines how many pixels are masked.

There are multiple significant differences between SGQN and MaDi. First of all, SGQN can be quite sensitive to the quantile hyperparameter. MaDi is free from this hyperparameter tuning, as it automatically finds the right fraction of pixels to mask. Furthermore, SGQN needs to compute gradients with respect to the inputs and weights, while MaDi only requires gradients of the weights. The additional components of SGQN are heavier, as they add about 1.6M parameters (an extra 25%) to the base architecture, with MaDi just adding roughly 10K (0.2%), reducing the memory requirements. MaDi also does not introduce any additional auxiliary loss function, as it is able to learn directly from the critic's objective. In essence, SGQN does not apply a mask to every input image, but uses them to learn better representations. MaDi tries to learn the most helpful masks such that the actor and critic receive only task-relevant information and are able to focus on the RL problem.

## 3 PRELIMINARIES

*Problem formulation.* We consider the problem of learning a policy for a Markov decision process (MDP) with the presence of visual distractions, similar to the formulation by Hansen et al. [21]. Our approach, MaDi, aims to learn a policy that generalizes well across MDPs with varying state spaces.

We formulate the interaction between the environment and policy as an MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$, where $\mathcal{S}$ is the state space, $\mathcal{A}$ is the action space, $\mathcal{P}: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the state transition function, $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\gamma$ is the discount factor. To address the challenges of partial observability [25], we define a state $\mathbf{s}_t$ as a sequence of $k$ consecutive frames $(\mathbf{o}_t, \mathbf{o}_{t-1}, \dots, \mathbf{o}_{t-(k-1)})$, $\mathbf{o}_i \in O$, where $O$ is the high-dimensional image space. In the particular benchmarks we employ for evaluation, $O = \mathbb{R}^{84 \times 84 \times 3}$ for the simulation environments and $O = \mathbb{R}^{160 \times 90 \times 3}$ for the robotic environment, as we receive RGB colored images as input with $84 \times 84$ and $160 \times 90$ pixels respectively.

Our goal is to learn a stochastic policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, where $\Delta(\mathcal{A})$ denotes the space of probability distributions over the action space $\mathcal{A}$. This policy aims to maximize the discounted return $R_t = \mathbb{E}_{\Gamma \sim \pi, \mathcal{P}} \left[ \sum_{t=0}^{T} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$ along a trajectory $\Gamma = (\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T)$. The policy $\pi$ is parameterized by a collection of learnable parameters $\theta$. We aim to learn parameters $\theta$ such that $\pi_\theta$ generalizes well across MDPs with perturbed observation spaces, denoted as $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$, where states $\overline{\mathbf{s}}_t \in \overline{\mathcal{S}}$ are constructed from observations $\overline{\mathbf{o}}_t \in \overline{O}$. The original observation space $O$ is a subset of the perturbed observation space $\overline{O}$, which may contain distractions.

*Distractions.* We define a distraction to be any input feature that is irrelevant to the task of the MDP, meaning that an optimal policy $\pi^*$ and value function $Q^*$ remain invariant under alterations of the feature value. In our case, input features are pixels $p_i \in \mathbb{R}^3$, where a state $\mathbf{s}$ consists of $n$ pixels: $\mathbf{s} = (p_1, p_2, \dots, p_n)$.

Suppose we encounter a state $\hat{\mathbf{s}}$, and we wish to determine whether pixel $p_i$ is a distraction in that particular state. Let $\mathcal{S}_i(\hat{\mathbf{s}})$ be the set of states where only pixel $p_i$ is changed in comparison to $\hat{\mathbf{s}}$. Then pixel $p_i$ is considered a distraction in state $\hat{\mathbf{s}}$ if $\pi^*$ and $Q^*$ remain invariant across the entire set $\mathcal{S}_i(\hat{\mathbf{s}})$. More formally:

DEFINITION 3.1. *A pixel $p_i$ is a **distraction** in state $\hat{\mathbf{s}}$ if, for an optimal policy $\pi^*$ and value function $Q^*$ it holds that, for an arbitrary but fixed action $\mathbf{a}$, we have $\forall \, \mathbf{s} \in \mathcal{S}_i(\hat{\mathbf{s}})$ :*

$$\pi^*(\mathbf{a}|\mathbf{s}) = \rho^* \qquad \text{where probability } \rho^* \in \mathbb{R} \text{ is constant,}$$
$$Q^*(\mathbf{s}, \mathbf{a}) = q^* \qquad \text{where value } q^* \in \mathbb{R} \text{ is constant.}$$

In other words: pixel $p_i$ can take on any value, but the optimal policy will not change. In that particular state $\hat{\mathbf{s}}$, the pixel $p_i$ is irrelevant to the task and thus a distraction. From this definition we can derive that the partial derivative of $\pi^*$ and $Q^*$ with respect to the input feature $p_i$ is zero.

COROLLARY 3.2. *If $p_i$ is a distraction in state $\hat{\mathbf{s}}$, then for an arbitrary action $\mathbf{a}$ we have that*

$$\frac{\partial}{\partial p_i} \pi^*(\mathbf{a}|\hat{\mathbf{s}}) = 0 \quad \text{and} \quad \frac{\partial}{\partial p_i} Q^*(\hat{\mathbf{s}}, \mathbf{a}) = 0.$$

This follows from Definition 3.1 since $\pi^*$ and $Q^*$ remain constant for varying $p_i$. Optimal policies perfectly ignore distractions, while suboptimal policies (i.e., neural networks during training) may be hindered by distractions. As distractions have no effect on the optimal policy, they can be safely masked when using $\pi^*$. This suggests that when striving to approximate $\pi^*$, it may be advantageous to mask distractions as well, a concept at the core of MaDi.

*Soft Actor-Critic.* In this work, we build upon the model-free off-policy reinforcement learning algorithm Soft Actor-Critic (SAC; [18]). SAC aims to estimate the optimal state-action value function $Q^*$ with its parameterized critic $Q_{\theta_Q}$. The actor is represented by a stochastic policy $\pi_{\theta_\pi}$, which aims to maximize the value outputted by the critic while simultaneously maintaining high entropy. The optional shared encoder $f_{\theta_f}$ is often used for SAC in image-based environments. The critic and shared encoder have target networks that start with the same parameters $\theta^{\text{tgt}} = \theta$. These are gradually updated throughout training by an exponential moving average: $\theta^{\text{tgt}} \longleftarrow (1 - \tau)\theta^{\text{tgt}} + \tau\theta$. We will often omit the implied parameter $\theta_N$ in our notation of any network $N$.

## 4 MADI

MaDi aims to mask distractions that hinder the agent from learning and performing well. We supplement the conventional actor-critic architecture of deep reinforcement learning agents by integrating a third, lightweight component, the Masker network $M$. The Masker adjusts the input by dimming irrelevant pixels, allowing the actor and critic networks to focus on learning the task at hand. The Masker and encoder compute internal representations using the Hadamard product and one forward call each: $f(\mathbf{s}_t \odot M(\mathbf{s}_t))$. Algorithm 1 indicates the few adjustments necessary to standard SAC (or SVEA, when using the augmentation on line 7). Note that MaDi does not need a target network for the Masker, reducing the additional number of parameters required.

**Algorithm 1** MaDi based on SAC

---

randomly initialize all networks: $\pi, Q, f, M$
copy parameters to target networks: $Q^{\text{tgt}}, f^{\text{tgt}}$

1: **for** timestep $t = 1...T$ **do**
    **act:**
2:     $\mathbf{a}_t \sim \pi\left(\cdot \,|\, f(\mathbf{s}_t \odot M(\mathbf{s}_t))\right)$               Sample action
3:     $\mathbf{s}_t', r_t \sim \mathcal{P}(\cdot \,|\, \mathbf{s}_t, \mathbf{a}_t)$          Perform action in env
4:     $\mathcal{B} \leftarrow \mathcal{B} \cup (\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_t')$          Add to replay buffer
    **update:**
5:     $\{\mathbf{s}_b, \mathbf{a}_b, r_b, \mathbf{s}_b'\} \sim \mathcal{B}$          Sample batch $b \subset \mathcal{B}$
6:     $\theta_\pi \leftarrow \theta_\pi - \eta \nabla_{\theta_\pi} \mathcal{L}_\pi(\mathbf{s}_b)$          Update $\pi$
7:     $\mathbf{s}_b \leftarrow \text{concat}(\mathbf{s}_b, \delta(\mathbf{s}_b))$          Apply augmentation
8:     **for** network $N$ in $[Q, f, M]$ **do**
9:         $\theta_N \leftarrow \theta_N - \eta \nabla_{\theta_N} \mathcal{L}_Q\left(\mathbf{s}_b, \mathbf{a}_b, r_b, \mathbf{s}_b'\right)$      Update $Q, f, M$
10:     **for** network $N$ in $[Q, f]$ **do**
11:         $\theta_N^{\text{tgt}} \leftarrow (1 - \tau)\theta_N^{\text{tgt}} + \tau \theta_N$          Update $Q^{\text{tgt}}, f^{\text{tgt}}$

---

## 4.1 The MaDi architecture

As shown in Figure 1, the Masker network is placed at the front of the agent's architecture. It produces a scalar multiplier for each pixel in the input image to determine the degree to which the pixel ought to be darkened. We refer to the output of this network as a *soft* mask.[3] These soft masks are applied element-wise to the observation frames, effectively reducing distractions (see Figure 2).

The Masker network is composed of three convolutional layers with ReLU non-linearities in between. The last layer outputs one channel with a Sigmoid activation function to squeeze values into the interval $[0, 1]$. The Masker receives three input channels, representing the RGB values of one frame $\boldsymbol{o}_t$. In Section 3, we defined a full state $\boldsymbol{s}_t$ to be a stack of $k$ frames, which is indeed what the actor and critic receive as input. The Masker is the only network that processes each frame separately. However, we still require only one forward pass through the Masker network for each input $\boldsymbol{s}_t$ of $k$ frames, as we efficiently reshape the channels into the batch dimension. See Appendix A for further implementation details.

## 4.2 How does the Masker learn?

One may expect that learning to output useful masks requires us to define a separate loss function, but this is not the case. The Masker can simply be updated via the critic's objective function[4] to update its parameters, as shown on line 9 of Algorithm 1. This means the masks are trained without any additional segmentation labels or saliency metrics. Our hypothesis on the surprising ability of MaDi to determine the task-relevant pixels solely from a scalar reward signal, pertains to the following:

- For *relevant* pixels, if the Masker network masks away essential pixels needed to determine an accurate Q-value, then the critic loss will presumably be high, and the Masker will thus be encouraged to leave these pixels visible.

---

[3] We also tried *hard* (i.e., *binary*) masks, but they proved more challenging to train.
[4] Future work could study whether the Masker network can also learn from the actor loss. In many SAC-based implementations with a shared ConvNet, the encoder is only updated by the critic loss, and it made sense to use these gradients for the Masker.

- For *irrelevant* pixels, we believe (and empirically show in Section 5.3) that strong and varying data augmentation helps. It gives an irrelevant pixel in a particular state $\boldsymbol{s}$ a varying pixel-value each time state $\boldsymbol{s}$ is sampled from the replay buffer, while the pixel's contribution to the $Q^*$-value remains the same (none, because it is irrelevant). The Masker is thus incentivized to mask this pixel, such that the actor and critic networks always see the same pixel-value for state $\boldsymbol{s}$, no matter which augmentation is used.

The Masker is updated together with the critic, which happens once for every environment step in case of synchronous runs. The robotic experiments use an asynchronous version of each algorithm. In that case, the Masker still gets as many updates as the critic, but it is no longer equal to the number of environment steps.

## 5 SIMULATION EXPERIMENTS

We present the experiments done on generalization benchmarks based on the DeepMind Control Suite [41] in this section, while our robotic experiments are shown in Section 6. We describe our experimental setup and provide the results obtained from our method, MaDi, compared with state-of-the-art approaches. Our experiments are designed to demonstrate the effectiveness of MaDi in masking distractions and improving generalization in vision-based RL.

## 5.1 Experimental Setup

*Benchmarks.* We evaluate the performance of MaDi on the Deep-Mind Control Generalization Benchmark [22, DMControl-GB] and the Distracting Control Suite [39, DistractingCS]. These benchmarks consist of a range of environments with varying levels of complexity and noise, providing a comprehensive assessment of an agent's ability to generalize to unseen, distracting environments.

- **DMControl-GB** has two setups with task-irrelevant pixels in the background: `video_easy` and `video_hard`. For the easy setup, there are just 10 videos to randomly sample from, and the surface from the training environment is still shown. In the hard counterpart, the surface is no longer visible, and one of 100 videos is selected. Note that all the frames in these videos are unseen — they do not overlap with the images in the augmentation dataset we use.
- **DistractingCS** goes a step further, also adjusting the camera's orientation and the agent's color during an episode, while displaying a randomly selected video in the background. The surface remains visible, such that the agent can orient itself during a changing camera angle. The intensity of the DistractingCS determines the difficulty. With higher intensity, the environment (1) samples from a larger set of videos, (2) changes the agent's color faster and to more extreme limits, and (3) adjusts the camera's orientation faster and to more severe angles. We use the default intensity level of 0.1. See the Supplementary Material for example videos.

*Environments.* Within these two benchmarks, we run all algorithms on six distinct environments, listed in Table 1. From the `cartpole` and `walker` domains we select two tasks, which differ by their starting positions and reward functions. See Appendix F for a detailed description of each task.

(a) `training_env (walker)`  (b) `video_hard (ball_in_cup)`  (c) `distracting_cs (cartpole)`

**Figure 2: Examples of original observations (left) and their masked versions (right) generated by MaDi in training (a) and testing (b, c) environments. Masks from other benchmark and domain combinations are shown in Appendix C.**
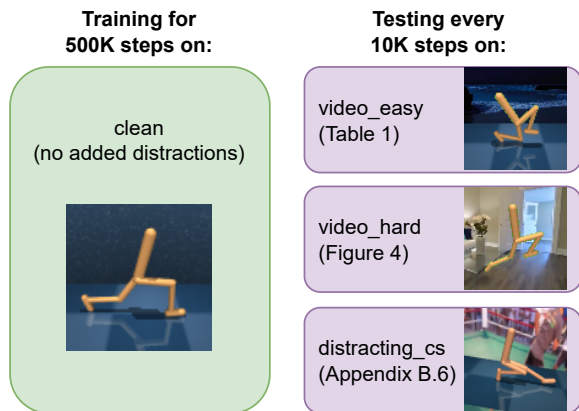


**Figure 3: The training and testing setup.**

**Table 1: Generalization performance of MaDi and various baseline algorithms on six different environments trained for** $500K$ **steps. We show undiscounted return on** `video_easy` **with mean and standard error over five seeds. MaDi outperforms or comes close to the state-of-the-art in all environments.**

| video_easy | SAC | DrQ | RAD | SODA | SVEA | SGQN | MaDi |
|---|---|---|---|---|---|---|---|
| ball_in_cup catch | 602 ±91 | 714 ±131 | 561 ±147 | 750 ±98 | 757 ±138 | 761 ±171 | **807** **±144** |
| cartpole balance | 924 ±19 | 932 ±33 | 801 ±95 | 961 ±10 | 967 ±2 | 965 ±5 | **982** **±4** |
| cartpole swingup | 782 ±21 | 613 ±74 | 658 ±17 | 215 ±125 | 786 ±15 | 798 ±13 | **848** **±6** |
| finger spin | 227 ±26 | 543 ±50 | 479 ±65 | 429 ±100 | 645 ±39 | 592 ±11 | **679** **±17** |
| walker stand | 507 ±113 | 954 ±10 | 961 ±1 | 147 ±17 | **977** **±3** | 672 ±153 | 967 ±3 |
| walker walk | 334 ±37 | 821 ±38 | 726 ±42 | 479 ±168 | **936** **±14** | 882 ±26 | 895 ±24 |
| avg | 563 | 763 | 698 | 497 | 845 | 778 | **863** |

*Models & Training.* For a fair comparison, we use the same base actor-critic architecture for all the methods considered in this study, including MaDi. All algorithms are trained for 500K timesteps on the clean training environment without distractions. We use the default hyperparameters for all baselines, as specified by the DMControl-GB [22]. See Appendix A for an overview of the hyperparameters.

*Augmentation.* As discussed in Section 2, all of our baselines (except SAC) use some form of data augmentation. MaDi is built on top of SVEA [21], which performs best with the `overlay` augmentation for distracting video backgrounds. Therefore, we choose to apply `overlay` for MaDi as well. This strong augmentation combines an observation frame from the training environment, $\mathbf{o}_t$, with a random image $\boldsymbol{x}$ from a large dataset as follows:

$$\delta_{\boldsymbol{x}}(\mathbf{o}_t) = \alpha \cdot \mathbf{o}_t + (1 - \alpha) \cdot \boldsymbol{x}$$

where $\delta$ denotes the augmentation function. In the experiments we use the default overlay factor of $\alpha = 0.5$ and sample images from the same dataset as used in SVEA, namely Places365 [59].

*Evaluation.* To assess generalization, we evaluate the trained agents zero-shot on a set of unseen environments with different levels of distractions. Specifically, we test on `video_easy` and `video_hard` from DMControl-GB, and on the Distracting Control Suite. Every 10K steps we evaluate the current policy for 20 episodes on the test environments; see Figure 3. We report the average undiscounted return over five random seeds during the last 10% of training, a metric often used to reduce variance [15, 16]. We run statistical tests to verify significance, shown in Appendix B.

## 5.2 Generalization Results

In Table 1 we show the results of MaDi and its baselines when generalizing to the `video_easy` setup of the DMControl-GB benchmark. MaDi is able to achieve the best or competitive performance in all environments. The learning curves of Figure 4 show that MaDi also generalizes well to the more challenging `video_hard` environments. Furthermore, the curves show that MaDi has a high sample efficiency, often reaching adequate performance in just 100K environment steps, by its ability to focus on the task-relevant pixels.

We present tables with results on the DistractingCS benchmark and the original training environments in Appendix B. MaDi also shows competitive performance in these additional settings. Note that the environments `ball_in_cup-catch` and `finger-spin` have *sparse* rewards, but even in this setting MaDi is able to perform well and generate useful masks. In Appendix C we provide examples of different masks produced for each domain.

## 5.3 Ablation on Augmentation

In Section 4.2 we described the expectation that MaDi would perform better with augmentations, as that can help it to recognize which pixels are irrelevant. To verify whether this intuition holds, we run five seeds of MaDi without the `overlay` augmentation on all
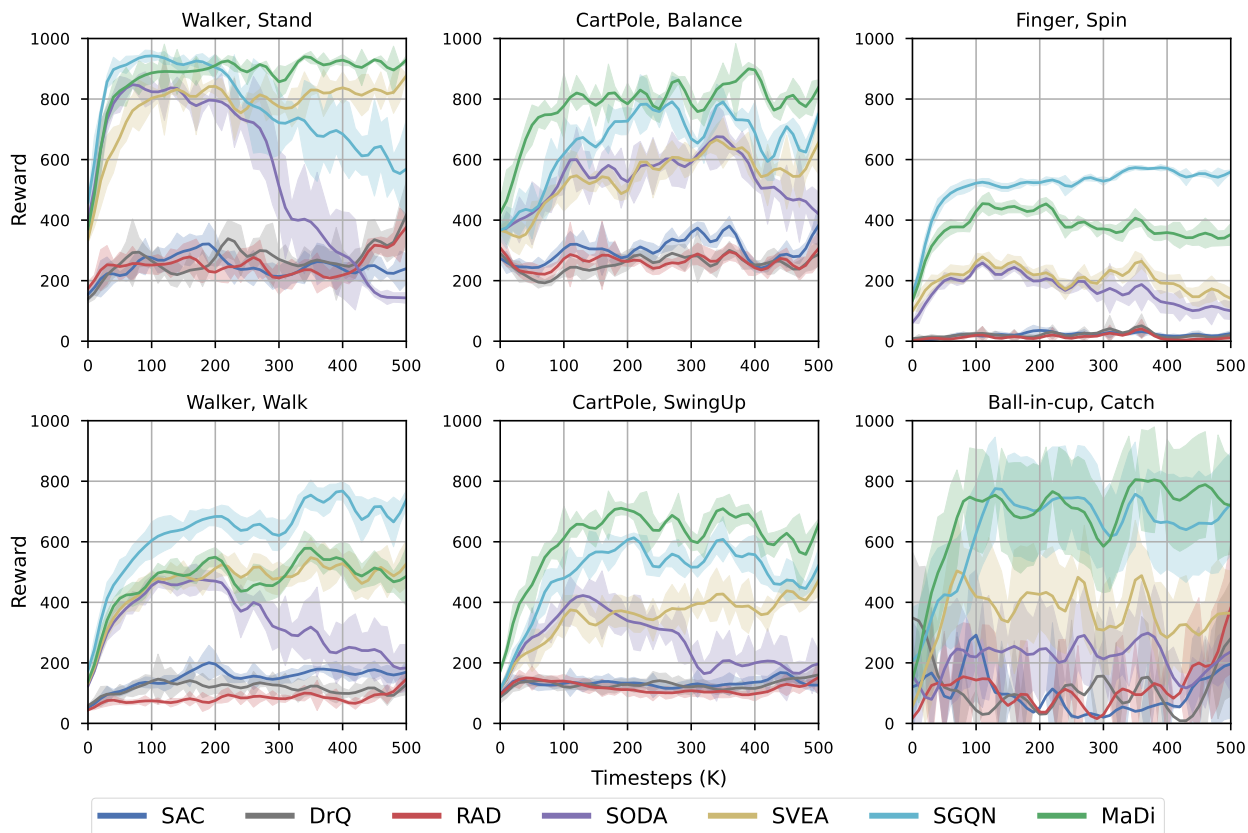
**Figure 4: Learning curves of MaDi and six baselines on `video_hard`. Agents are trained on clean data for $500K$ steps and tested on `video_hard` every $10K$ steps. MaDi often reaches the top of the class, while some baselines can overfit to the training environment and decrease in generalizability. The curves show the mean over five seeds with standard error shaded alongside.**

**Table 2: Ablation study showing the effect of the `overlay` augmentation. SVEA and MaDi both use it, while SAC and MaDi-SAC do not. All algorithms are trained for $500K$ steps. We show mean undiscounted return and standard error over five seeds evaluated on `video_hard`. The results reveal that MaDi benefits from data augmentation.**

| video_hard | SAC | MaDi-SAC | SVEA | MaDi |
|---|---|---|---|---|
| ball_in_cup catch | 176 ±38 | 190 ±52 | 327 ±59 | 758 ±135 |
| cartpole balance | 314 ±12 | 237 ±6 | 579 ±26 | 827 ±25 |
| cartpole swingup | 140 ±10 | 132 ±9 | 453 ±26 | 619 ±24 |
| finger spin | 21 ±4 | 121 ±36 | 154 ±31 | 358 ±25 |
| walker stand | 233 ±28 | 320 ±88 | 847 ±18 | 920 ±14 |
| walker walk | 168 ±19 | 95 ±24 | 526 ±55 | 504 ±33 |
| avg | 175 | 183 | 481 | 664 |

six environments. We call this variant MaDi-SAC, as it now builds on top of SAC [18] instead of SVEA.

The results are shown in Table 2. The generalization performance of the algorithms that use augmentation is much better than those without. MaDi indeed benefits from data augmentation to become more robust against previously unseen distractions. Furthermore, in Appendix D we present results on a few other types of augmentations that could be combined with our method.

## 5.4 Does MaDi work with Vision Transformers?

In image-based RL, the use of Vision Transformers [8, **ViT**] has recently gained in popularity [21, 40]. We set up a small experiment to verify whether MaDi still works in this setting. The shared encoder (see Figure 1), which is originally an 11-layer ConvNet, is now replaced by a ViT of 4 blocks with 8 attention heads each. We maintain the same architecture for the Masker network. More implementation details on the ViT encoder are in Appendix A.1.

We train the ViT-based versions of SVEA and MaDi for 300K timesteps on the clean environments of `walker-walk` and `cartpole-swingup`, and test on `video_hard` at every 10K steps. The results presented in Figure 5 show that not only does MaDi work well with a ViT encoder, but it even boosts the generalization performance.
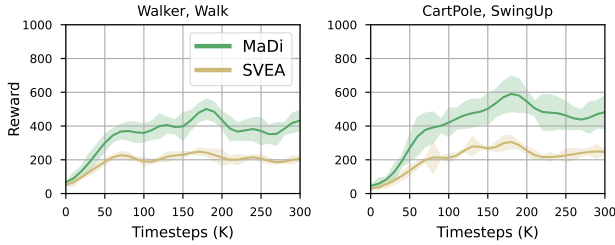
**Figure 5: Generalization performance of MaDi and SVEA on `video_hard` when trained with a ViT encoder for** $300K$ **steps. We show the mean and standard error over five seeds.**

# 6 ROBOTIC EXPERIMENTS

In this section, we describe our experiments on a robotic arm visual reacher task, showing that MaDi can learn to generalize when facing distracting backgrounds, even in real-world environments.

## 6.1 Experimental Setup

*Robotic Arm.* The UR5 industrial robot arm consists of six joints that can rotate to move the tip to a desired position. It uses the conventional TCP/IP protocol to transmit the state of the arm to the host computer and receive actuation packets in return at an interval of 8 milliseconds. A state packet includes the angles, velocities, target accelerations, and currents of all six joints. We configured the UR5 to use the velocity control mode. Since the UR5 does not come with a camera, we attached a Logitech RGB camera to the tip of the arm to facilitate vision-based tasks. We use *SenseAct*, a computational framework for robotic learning experiments to communicate with the robot [30]. This setup enables robust and reproducible learning across different locations and conditions.

*Training environment.* We train on the UR5-VisualReacher task [48, 52] for real-world robot experiments. Figure 6 depicts the experimental setup. This task involves using a camera to guide a UR5 robotic arm to reach a red target on a monitor. We ensure that the arm stays within a safe bounding box to avoid collisions. The available actions are represented by the desired angular velocities for five[5] joints, ranging from −0.7 to 0.7 radians per second. The full observation that the learning agent receives includes three consecutive RGB images from the Logitech camera of dimensions $160 \times 90$ pixels, joint angles, joint angular velocities, and the previous action taken. The reward function is defined as follows [52]:

$$r_t = \frac{c}{hw} M_t \odot W$$

where $c$ is a scaling coefficient, $h$ and $w$ are the height and width of the image in pixels, $M_t$ is a binary mask[6] of shape $h \times w$ that detects for each pixel whether it is currently red, and $W$ is a weighting matrix of shape $h \times w$ with values decreasing from 1 at its center to 0 near the edges. These are multiplied element-wise by the Hadamard product $\odot$. The reward incentivizes the robot to move its camera closer to the target and keep the target at the center of the frame. We set the coefficient $c = 800$ for all experiments and clip the rewards to be between 0 and 4. An episode lasts 150 timesteps of 40 ms each,

---

[5]We do not move the sixth joint because its sole purpose is to control a gripper.
[6]Note that this is a preprogrammed mask based on RGB thresholds, not made by MaDi.



| (a) Training environment | (b) Testing environment |

**Figure 6: The UR5-VisualReacher(-VideoBackgrounds) benchmark used in our experiments. The agent is rewarded for getting its camera as close as possible to the red circle randomly located on the screen. For visual demonstrations of the robot's performance with each algorithm, please refer to the videos in the Supplementary Material.**

for 6 seconds in total. The agent sends an action at every timestep, which is repeated by the SenceAct system five times at every 8 ms.

*Real-time asynchronous learning.* We use the ReLoD system [48] to facilitate effective, real-time learning. In simulation, environments can be internally paused while agents carry out learning updates. However, in the real world, the environment does not wait for the agent to complete its sequential computations and learning updates. We use an asynchronous implementation of SAC [48, 52] to gather data in real-time and implement learning updates through a separate process. We reimplemented MaDi and all baselines used in these experiments to their asynchronous versions. Each learning update can take anywhere from 70 to 500 ms, depending on the chosen algorithm and hyperparameter configuration, meaning that there will be fewer updates than steps in the environment.

*Testing environment.* We test the generalization performance on a similar task, but now with videos playing in the background on the screen. We define this new generalization benchmark as UR5-VisualReacher-VideoBackgrounds (see Figure 6b), selecting five videos from the DMControl-GB [22] to use as backgrounds. We test on this benchmark after every 4500 timesteps in the training environment. This evaluation is always done for 10 episodes, twice on each video. See Appendix F for examples of frames.

*Baselines.* We compare the performance of SAC [18], RAD [28], SVEA [21], and MaDi. The base architecture is the same for all algorithms, but it differs somewhat from the simulation experiments of Section 5. See Appendix A for more details.

## 6.2 Results

In Figure 7, we present the results on the testing environment with video backgrounds. Without being trained on these distracting visuals, MaDi is able to generalize well in this challenging task.

The algorithms are trained asynchronously in this robotic environment, which means that MaDi will make fewer updates as it uses the additional Masker network. However, as Figure 7 shows, this only incurs a small delay in learning while gaining superior generalization capability through the increased focus on task-relevant pixels. A similar pattern is present on the original training environment, as shown in Figure 10 of Appendix B.1.
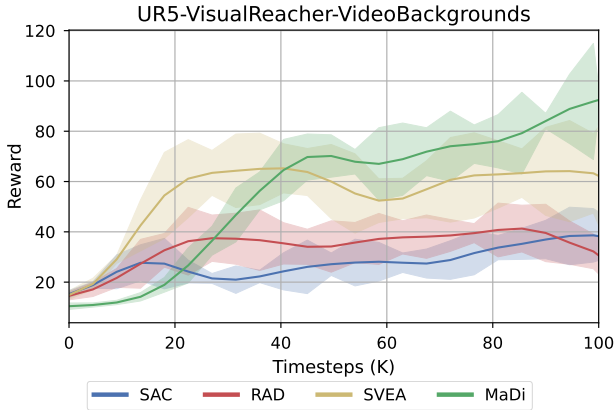
Figure 7: **Performance of MaDi and three baselines on the UR5-VisualReacher task with random videos playing in the background. Agents are trained on the clean environment for** $100K$ **steps. We display the mean and standard error over five seeds. MaDi already generalizes best after** $50K$ **timesteps.**

Table 3: **The average reward per timestep on the UR5-VisualReacher task during the last 10% of steps in an episode. MaDi receives higher rewards in the final position of an episode in both training and testing environments, showing that it finds the red target with higher accuracy.**

| Reward per step | SAC | RAD | SVEA | MaDi |
|---|---|---|---|---|
| Training env. | 1.38 ±0.09 | 1.32 ±0.11 | 1.18 ±0.37 | 1.95 ±0.09 |
| Testing env. | 0.32 ±0.04 | 0.24 ±0.07 | 0.47 ±0.14 | 0.74 ±0.07 |

In these experiments, both MaDi and SVEA use the `overlay` augmentation. See Appendix D.1 for results showing the effect of the `conv` augmentation in this environment. MaDi also outperforms the baselines with the `conv` augmentation but scores slightly lower, despite the arguably clearer masks it generates in that setting.

In Appendix B.2, we present an additional experiment with a *sparse reward* function. In this UR5-VisualReacher-SparseRewards task, the agent is only rewarded with a +1 whenever its camera is close enough to the red target (according to a predefined threshold), and receives zero reward otherwise. Even in this challenging setting, MaDi is able to surpass the baselines in generalization performance.

## 6.3 Analysis

In Figure 8 we show that MaDi can learn to recognize the task-relevant features even in this real-world robotic task. It is able to generalize to the unseen testing environment and produce helpful masks. There seem to be fewer pixels dimmed in this environment compared to the other benchmarks, which may be because it can be useful for the agent to know where the entire screen is positioned.

*Distinction between training and testing.* MaDi performs well in both the training and test environments, but there is quite a large discrepancy between the total rewards received. For all algorithms, the reward in the testing environment is substantially lower than in
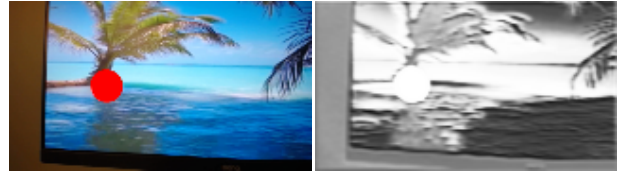


Figure 8: **Observation and the corresponding mask generated by MaDi in the UR5-VisualReacher-VideoBackgrounds task. The mask is subtle, but clearly leaves the red dot's pixels intact while dimming other areas of the frame. See Appendix C for more examples of masks.**

the training environment. Taking a qualitative look at the behavior of the robotic arm, it seems this is mostly due to the fact that the robotic arm moves slower toward the target in the testing environment than in the training environment. The agents encounter unseen observations that significantly deviate from the training environment, likely driving them to select different actions that do not match well in sequence, causing the arm to slow down. The SAC and RAD agents rarely complete the task at all when there are videos playing in the background.

Even though the movement towards the goal is slower in the testing environment for all algorithms, MaDi does often reach a (near) optimal state at the end of its trajectory, similar to training. See Table 3 for an overview of the rewards in the last 10% of steps for each baseline. MaDi shows a higher accuracy in finding the red target near the end of an episode.[7]

## 7 CONCLUSION

In the domain of vision-based deep reinforcement learning, we formalize the problem setting of distracting task-irrelevant features. We propose a novel method, MaDi, which learns directly from the reward signal to *ma*sk *di*stractions with a lightweight Masker network, without requiring any additional segmentation labels or loss functions. Our experiments show that MaDi is competitive with state-of-the-art algorithms on the DeepMind Control Generalization Benchmark and the Distracting Control Suite, while only using 0.2% additional parameters. The masks generated by MaDi enhance the agent's focus by dimming visual distractions. Even in the sparse reward setting, the Masker network is able to learn where the task-relevant pixels are in each state. Furthermore, we test MaDi on a real UR5 Robotic Arm, showing that it can outperform the baselines not only in simulation environments but also on our newly defined UR5-VisualReacher-VideoBackgrounds generalization benchmark.

*Limitations & Future Work.* The algorithms in this work build on the model-free off-policy deep RL algorithm SAC, while other options remain open for investigation. In future work, we seek to apply MaDi to other reinforcement learning algorithms such as PPO [38] or DQN [32] and improve the clarity of its masks. We have experimented with MaDi on one robotic arm, it would be interesting to see whether the Masker network can also produce useful masks on a diverse set of robots. Lastly, in future research, we aim to explore the possibility of using MaDi for transfer learning.

---

[7] Moreover, MaDi generally moves to the target faster, see youtu.be/TQMazg6dntE.

## ETHICS STATEMENT

Our research aims to contribute to the development of reinforcement learning algorithms to facilitate their application in practical scenarios that positively impact society. We believe our work with MaDi has the potential to contribute to, for instance, enhancing a household robot's ability to focus on relevant visual information amidst a clutter of distractions. Furthermore, we aim to minimize the computational footprint of our algorithms by adding a negligible amount of parameters to the original structure, supporting the development of sustainable and energy-efficient AI.

## REFERENCES

[1] Rishabh Agarwal, Marlos C Machado, Pablo Samuel Castro, and Marc G Bellemare. 2020. Contrastive Behavioral Similarity Embeddings for Generalization in Reinforcement Learning. In *International Conference on Learning Representations*. URL: https://arxiv.org/abs/2101.05265. (Cited in Section 2)

[2] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. 2019. Solving Rubik's Cube with a Robot Hand. *arXiv preprint arXiv:1910.07113* (2019). URL: https://openai.com/research/solving-rubiks-cube. (Cited in Section 1, 2)

[3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer Normalization. *Advances in Neural Information Processing Systems, Deep Learning Symposium* (2016). URL: https://arxiv.org/abs/1607.06450. (Cited in Section A)

[4] David Bertoin, Adil Zouitine, Mehdi Zouitine, and Emmanuel Rachelson. 2022. Look where you look! Saliency-guided Q-networks for generalization in visual Reinforcement Learning. *Advances in Neural Information Processing Systems* 35 (2022), 30693–30706. URL: https://arxiv.org/abs/2209.09203. (Cited in Section 1, 2, 2, 2, 4)

[5] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym. *arXiv preprint arXiv:1606.01540* (2016). URL: https://www.gymlibrary.dev/. (Cited in Section 2)

[6] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. 2019. Quantifying Generalization in Reinforcement Learning. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 1282–1289. URL: https://proceedings.mlr.press/v97/cobbe19a.html. (Cited in Section 2)

[7] Jonas Degrave, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, et al. 2022. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature* 602, 7897 (2022), 414–419. URL: https://www.nature.com/articles/s41586-021-04301-9. (Cited in Section 1)

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations* (2021). URL: https://arxiv.org/abs/2010.11929. (Cited in Section 5.4, A.1)

[9] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. 2016. RL$^2$: Fast Reinforcement Learning via Slow Reinforcement Learning. *arXiv preprint arXiv:1611.02779* (2016). URL: https://arxiv.org/abs/1611.02779. (Cited in Section 2)

[10] Open X-Embodiment Collaboration et al. 2023. Open X-Embodiment: Robotic Learning Datasets and RT-X Models. URL: https://robotics-transformer-x.github.io. (Cited in Section 1)

[11] Linxi Fan, Guanzhi Wang, De-An Huang, Zhiding Yu, Li Fei-Fei, Yuke Zhu, and Anima Anandkumar. 2021. SECANT: Self-Expert Cloning for Zero-Shot Generalization of Visual Policies. *International Conference on Machine Learning* (2021). URL: https://arxiv.org/abs/2106.09678. (Cited in Section 2)

[12] Jesse Farebrother, Marlos C Machado, and Michael Bowling. 2018. Generalization and Regularization in DQN. *NeurIPS'18 Deep Reinforcement Learning Workshop* (2018). URL: https://arxiv.org/abs/1810.00123. (Cited in Section 2)

[13] Norm Ferns, Prakash Panangaden, and Doina Precup. 2011. Bisimulation Metrics for Continuous Markov Decision Processes. *SIAM J. Comput.* 40, 6 (2011), 1662–1714. URL: https://doi.org/10.1137/10080484X. (Cited in Section 2)

[14] Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. 2016. Deep Spatial Autoencoders for Visuomotor Learning. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 512–519. URL: https://arxiv.org/abs/1509.06113. (Cited in Section A)

[15] Laura Graesser, Utku Evci, Erich Elsen, and Pablo Samuel Castro. 2022. The State of Sparse Training in Deep Reinforcement Learning. In *International Conference on Machine Learning*. PMLR, 7766–7792. URL: https://arxiv.org/abs/2206.10369. (Cited in Section 5.1)

[16] Bram Grooten, Ghada Sokar, Shibhansh Dohare, Elena Mocanu, Matthew E. Taylor, Mykola Pechenizkiy, and Decebal Constantin Mocanu. 2023. Automatic Noise Filtering with Dynamic Sparse Training in Deep Reinforcement Learning. *The 22nd International Conference on Autonomous Agents and Multiagent Systems (AAMAS)* (2023). URL: https://arxiv.org/abs/2302.06548. (Cited in Section 2, 5.1)

[17] Tuomas Haarnoja, Ben Moran, Guy Lever, Sandy H Huang, Dhruva Tirumala, Markus Wulfmeier, Jan Humplik, Saran Tunyasuvunakool, Noah Y Siegel, Roland Hafner, et al. 2023. Learning Agile Soccer Skills for a Bipedal Robot with Deep Reinforcement Learning. *arXiv preprint arXiv:2304.13653* (2023). URL: https://sites.google.com/view/op3-soccer. (Cited in Section 1)

[18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *International Conference on Machine Learning*. PMLR, 1861–1870. URL: https://arxiv.org/abs/1801.01290. (Cited in Section 1, 2, 3, 5.3, 6.1, A)

[19] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. 2020. Dream to Control: Learning Behaviors by Latent Imagination. *International Conference on Learning Representations* (2020). URL: https://arxiv.org/abs/1912.01603. (Cited in Section F.1)

[20] Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà, Pieter Abbeel, Alexei A Efros, Lerrel Pinto, and Xiaolong Wang. 2020. Self-Supervised Policy Adaptation during Deployment. *International Conference on Learning Representations* (2020). URL: https://arxiv.org/abs/2007.04309. (Cited in Section 2)

[21] Nicklas Hansen, Hao Su, and Xiaolong Wang. 2021. Stabilizing Deep Q-Learning with ConvNets and Vision Transformers under Data Augmentation. *35th Conference on Neural Information Processing Systems* (2021). URL: https://arxiv.org/abs/2107.00644. (Cited in Section 1, 2, 2, 3, 5.1, 5.4, 6.1, 4, A.1, F.1)

[22] Nicklas Hansen and Xiaolong Wang. 2021. Generalization in Reinforcement Learning by Soft Data Augmentation. In *International Conference on Robotics and Automation*. URL: https://arxiv.org/abs/2011.13389. (Cited in Section 1, 1, 2, 2, 2, 5.1, 6.1, A, 4, A, B.5, D, F.1, F.2)

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 770–778. URL: https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/He_Deep_Residual_Learning_CVPR_2016_paper.pdf. (Cited in Section 2)

[24] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian Error Linear Units (GELUs). *arXiv preprint arXiv:1606.08415* (2016). URL: https://arxiv.org/abs/1606.08415. (Cited in Section A.1)

[25] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. 1998. Planning and acting in partially observable stochastic domains. *Artificial Intelligence* 101, 1-2 (1998), 99–134. URL: https://www.sciencedirect.com/science/article/pii/S000437029800023X. (Cited in Section 3)

[26] Diederik Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. *International Conference for Learning Representations* (2015). URL: https://arxiv.org/abs/1412.6980. (Cited in Section 4)

[27] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. 2023. A Survey of Zero-shot Generalisation in Deep Reinforcement Learning. *Journal of Artificial Intelligence Research* 76 (2023), 201–264. URL: https://arxiv.org/abs/2111.09794. (Cited in Section 2)

[28] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. 2020. Reinforcement Learning with Augmented Data. *Advances in*

*Neural Information Processing Systems* 33 (2020), 19884–19895. URL: https://arxiv.org/abs/2004.14990. (Cited in Section 1, 2, 2, 2, 6.1, 4)

[29] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. 2020. CURL: Contrastive Unsupervised Representations for Reinforcement Learning. In *International Conference on Machine Learning*. PMLR, 5639–5650. URL: https://arxiv.org/abs/2004.04136. (Cited in Section 2)

[30] A Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra. 2018. Benchmarking Reinforcement Learning Algorithms on Real-World Robots. In *Conference on robot learning*. PMLR, 561–591. URL: https://arxiv.org/abs/1809.07731. (Cited in Section 6.1)

[31] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, et al. 2016. Learning to Navigate in Complex Environments. *arXiv preprint arXiv:1611.03673* (2016). URL: https://arxiv.org/abs/1611.03673. (Cited in Section 1)

[32] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533. URL: https://www.nature.com/articles/nature14236. (Cited in Section 1, 2, 7)

[33] Xinghua Qu, Zhu Sun, Yew-Soon Ong, Abhishek Gupta, and Pengfei Wei. 2020. Minimalistic Attacks: How Little it Takes to Fool a Deep Reinforcement Learning Policy. *IEEE Transactions on Cognitive and Developmental Systems* 13, 4 (2020), 806–817. URL: https://arxiv.org/abs/1911.03849. (Cited in Section 2)

[34] Roberta Raileanu, Maxwell Goldstein, Denis Yarats, Ilya Kostrikov, and Rob Fergus. 2021. Automatic Data Augmentation for Generalization in Reinforcement Learning. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (Eds.), Vol. 34. Curran Associates, Inc., 5402–5415. URL: https://proceedings.neurips.cc/paper/2021/hash/2b38c2df6a49b97f706ec9148ce48d86-Abstract.html. (Cited in Section 2)

[35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*. Springer, 234–241. URL: https://arxiv.org/abs/1505.04597. (Cited in Section 2)

[36] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* 115, 3 (2015), 211–252. https://doi.org/10.1007/s11263-015-0816-y URL: https://link.springer.com/article/10.1007/s11263-015-0816-y. (Cited in Section 2)

[37] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. 2020. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature* 588, 7839 (2020), 604–609. URL: https://www.nature.com/articles/s41586-020-03051-4. (Cited in Section 1)

[38] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347* (2017). URL: https://arxiv.org/abs/1707.06347. (Cited in Section 7)

[39] Austin Stone, Oscar Ramirez, Kurt Konolige, and Rico Jonschkowski. 2021. The Distracting Control Suite – A Challenging Benchmark for Reinforcement Learning from Pixels. *arXiv preprint arXiv:2101.02722* (2021). URL: https://arxiv.org/abs/2101.02722. (Cited in Section 1, 1, 2, 5.1, B.6, F.1)

[40] Tianxin Tao, Daniele Reda, and Michiel van de Panne. 2022. Evaluating Vision Transformer Methods for Deep Reinforcement Learning from Pixels. *arXiv preprint arXiv:2204.04905* (2022). URL: https://arxiv.org/abs/2204.04905. (Cited in Section 5.4)

[41] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. 2018. Deepmind Control Suite. *arXiv preprint arXiv:1801.00690* (2018). URL: https://arxiv.org/abs/1801.00690. (Cited in Section 2, 5, F.1)

[42] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. 2017. Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 23–30. URL: https://arxiv.org/abs/1703.06907. (Cited in Section 2)

[43] Emanuel Todorov, Tom Erez, and Yuval Tassa. 2012. MuJoCo: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 5026–5033. URL: https://mujoco.org/. (Cited in Section 2)

[44] Manan Tomar, Riashat Islam, Sergey Levine, and Philip Bachman. 2023. Ignorance is Bliss: Robust Control via Information Gating. *arXiv preprint arXiv:2303.06121* (2023). URL: https://arxiv.org/abs/2303.06121. (Cited in Section 2)

[45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *Advances in Neural Information Processing Systems* 30 (2017). URL: https://arxiv.org/abs/1706.03762. (Cited in Section A.1)

[46] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. 2016. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763* (2016). URL: https://arxiv.org/abs/1611.05763. (Cited in Section 2)

[47] Kaixin Wang, Bingyi Kang, Jie Shao, and Jiashi Feng. 2020. Improving Generalization in Reinforcement Learning with Mixture Regularization. *Advances in Neural Information Processing Systems* 33 (2020), 7968–7978. URL: https://arxiv.org/abs/2010.10814. (Cited in Section 2)

[48] Yan Wang, Gautham Vasan, and A Rupam Mahmood. 2023. Real-Time Reinforcement Learning for Vision-Based Robotics Utilizing Local and Remote Computers. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 9435–9441. URL: https://arxiv.org/abs/2210.02317. (Cited in Section 6.1, 6.1, A, F.2)

[49] Jinwei Xing, Takashi Nagata, Kexin Chen, Xinyun Zou, Emre Neftci, and Jeffrey L Krichmar. 2021. Domain Adaptation In Reinforcement Learning Via Latent Unified State Representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 10452–10459. URL: https://arxiv.org/abs/2102.05714. (Cited in Section 2)

[50] Denis Yarats, Ilya Kostrikov, and Rob Fergus. 2021. Image Augmentation Is All You Need: Regularizing Deep Reinforcement Learning from Pixels. *International Conference on Learning Representations* (2021). URL: https://openreview.net/forum?id=GY6-6sTvGaf. (Cited in Section 1, 2, 2, 2, 4, F.1)

[51] Tao Yu, Zhizheng Zhang, Cuiling Lan, Yan Lu, and Zhibo Chen. 2022. Mask-based Latent Reconstruction for Reinforcement Learning. *Advances in Neural Information Processing Systems* 35 (2022), 25117–25131. URL: https://arxiv.org/abs/2201.12096. (Cited in Section 2)

[52] Yufeng Yuan and A Rupam Mahmood. 2022. Asynchronous Reinforcement Learning for Real-Time Control of Physical Robots. In *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 5546–5552. URL: https://arxiv.org/abs/2203.12759. (Cited in Section 6.1, 6.1, F.2)

[53] Zhecheng Yuan, Guozheng Ma, Yao Mu, Bo Xia, Bo Yuan, Xueqian Wang, Ping Luo, and Huazhe Xu. 2022. Don't Touch What Matters: Task-Aware Lipschitz Data Augmentation for Visual Reinforcement Learning. *International Joint Conference on Artificial Intelligence* (2022). URL: https://arxiv.org/abs/2202.09982. (Cited in Section 2, 2)

[54] Zhecheng Yuan, Zhengrong Xue, Bo Yuan, Xueqian Wang, Yi Wu, Yang Gao, and Huazhe Xu. 2022. Pre-Trained Image Encoder for Generalizable Visual Reinforcement Learning. *Advances in Neural Information Processing Systems* 35 (2022), 13022–13037. URL: https://arxiv.org/abs/2212.08860. (Cited in Section 2, 2)

[55] Dylan Yung, Andrew Szot, Prithvijit Chattopadhyay, Judy Hoffman, and Zsolt Kira. 2023. Augmentation Curriculum Learning For Generalization in RL. (2023). URL: https://openreview.net/forum?id=Fj1S0SV8p3U. (Cited in Section D)

[56] Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. 2021. Learning Invariant Representations for Reinforcement Learning without Reconstruction. *International Conference on Learning Representations (ICLR)* (2021). URL: https://arxiv.org/abs/2006.10742. (Cited in Section 2)

[57] Chiyuan Zhang, Oriol Vinyals, Remi Munos, and Samy Bengio. 2018. A Study on Overfitting in Deep Reinforcement Learning. *arXiv preprint arXiv:1804.06893* (2018). URL: https://arxiv.org/abs/1804.06893. (Cited in Section 2)

[58] Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui Hsieh. 2020. Robust Deep Reinforcement Learning against Adversarial Perturbations on State Observations. *Advances in Neural Information Processing Systems* 33 (2020), 21024–21037. URL: https://arxiv.org/abs/2003.08938. (Cited in Section 2)

[59] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Places: A 10 million Image Database for Scene Recognition. *IEEE transactions on pattern analysis and machine intelligence* 40, 6 (2017), 1452–1464. URL: https://ieeexplore.ieee.org/document/7968387. (Cited in Section 2, 5.1)

[60] Y Zhu, R Mottaghi, E Kolve, JJ Lim, A Gupta, L Fei-Fei, and A Farhadi. 2017. Target-driven Visual Navigation in Indoor Scenes using Deep Reinforcement Learning. *International Conference on Robotics and Automation* (2017). URL: https://ieeexplore.ieee.org/document/7989381. (Cited in Section 1)

# APPENDIX

## A   IMPLEMENTATION DETAILS

We provide further details on our implementation of MaDi. Table 4 presents an overview of the hyperparameters used in our study. For the simulation experiments we used the default hyperparameters as specified by the DMControl Generalization Benchmark [22] and the particular baselines. For the robotic task, we obtained hyperparameters from [48] while aiming to keep the majority the same.

Table 4: Overview of the hyperparameters used in our experiments.

| Hyperparameter | Simulation experiments | Robotic experiments (if different) |
| --- | --- | --- |
| *Shared by all algorithms* | | |
| optimizer | Adam [26] | |
| learning rate (actor and critic) | $10^{-3}$ | $3 \cdot 10^{-4}$ |
| Adam $\beta_1, \beta_2$ (all networks) | 0.9, 0.999 | |
| discount ($\gamma$) | 0.99 | |
| frame stack | 3 | |
| action repeat | 4 | 5 |
| batch size | 128 | |
| initial collect steps | 1000 | |
| replay buffer size | 500K | 100K |
| environment steps | 500K | 100K |
| train every $k_c$ env steps (critic), $k_c =$ | 1 | asynchronous |
| train every $k_a$ env steps (actor), $k_a =$ | 2 | once every critic update |
| gradient steps per training step | 1 | |
| target update interval ($k_{tar}$) | 2 | once every critic update |
| target smoothing coefficient critic ($\tau$) | 0.01 | 0.005 |
| target smoothing coefficient encoder ($\tau_{\mathrm{enc}}$) | 0.05 | 0.005 |
| initial temperature ($\alpha$ in SAC) | 0.1 | |
| learning rate (for temperature $\alpha$) | $10^{-4}$ | $3 \cdot 10^{-4}$ |
| Adam $\beta_1, \beta_2$ (for temperature $\alpha$) | 0.5, 0.999 | |
| *DrQ* [50] | | |
| K, M ($Q$-target, $Q$-function averaging) | 1, 1 | |
| random shift | Up to ±4 pixels | |
| *RAD* [28] | | |
| random crop | $100 \times 100 \longrightarrow 84 \times 84$ | $160 \times 90 \longrightarrow 156 \times 88$ |
| *SODA* [22] | | |
| learning rate auxiliary network | $10^{-3}$ | |
| batch size auxiliary network | 256 | |
| target smoothing coefficient aux net ($\tau_{\mathrm{aux}}$) | 0.005 | |
| train every $k_{\mathrm{aux}}$ env steps (aux), $k_{\mathrm{aux}} =$ | 2 | |
| *SVEA* [21] | | |
| svea_alpha, svea_beta | 0.5, 0.5 | |
| *SGQN* [4] | | |
| threshold sgqn_quantile[8] | 0.95 or 0.98 | |
| weight decay critic | $10^{-5}$ | |
| learning rate SGQN head | $3 \cdot 10^{-4}$ | |
| train every $k_{\mathrm{sgqn}}$ env steps (SGQN head), $k_{\mathrm{sgqn}} =$ | 2 | |
| *MaDi* | | |
| learning rate masker network | $10^{-3}$ | $3 \cdot 10^{-4}$ |
| train every $k_{\mathrm{m}}$ env steps (masker), $k_{\mathrm{m}} =$ | 1 | |

---

[8]The sgqn_quantile is set to 0.95 for domains with large agents (e.g., walker and finger), and 0.98 for domains with smaller agents (e.g., cartpole and ball_in_cup).

*Network Architecture.* All the algorithms in our experiments use the same base architecture, acquired from the DMControl Generalization Benchmark [22], to ensure a fair comparison. This base architecture (yellow modules in Figure 1) consists of:

- **Encoder:** a shared 11-layer ConvNet that takes a stack of 3 RGB frames (shape $[9, 84, 84]$) and outputs spatial features of size $[32, 21, 21]$. Each layer has 32 channels with a kernel size of $3 \times 3$, a stride of 1 and no padding. Only the first layer has a stride of 2, to reduce the channel size quicker. We use ReLU activations between all convolutional layers.
- **Actor:** a multi-layer perceptron (MLP) that outputs actions. It first projects the $32 \cdot 21 \cdot 21 = 14,112$ flattened features to a vector of just 100 through a large MLP layer. We then apply LayerNorm [3] and a Tanh activation function. The bulk of the actor network follows, with three MLP layers of hidden dimension 1024 with ReLU activations in between. The output layer has twice the size of the action dimension, to output means $\mu$ and standard deviations $\sigma$.
- **Critic:** two independent MLPs that output state-action values. Both have the same architecture as the actor. The large projection (up to and including the Tanh activation) is shared between the two critics. The critics receive the current action $a$ as input as well, this is concatenated to the output of the large projection. Their last layers only have one output neuron, determining the $Q$-value estimate.

For the robotic experiments the architectures are a bit different:

- **Encoder:** a shared 4-layer ConvNet that takes a stack of 3 RGB frames (shape $[9, 90, 160]$), followed by a spatial softmax layer [14]. We use 32 channels in each layer, kernel size is $3 \times 3$, no padding, and stride 2 on all layers except the last one, which uses 1. ReLU activation functions between the layers. The encoder outputs a vector of length 64.
- **Actor:** an MLP that outputs actions. It first concatenates the encoder features with a proprioception vector of length 15, which consists of the angles, angular velocities, and previous actions of each of the five joints that we control. (The sixth joint is for a gripper, which we do not use.) After that, two MLP layers of hidden dimension 512 follow.
- **Critic:** two independent MLPs that output state-action values. Both have the same architecture as the actor, except again the concatenation of the action $a$ in the input and an output layer of dimension 1.

And of course, the critic network has an exact replica, its *target network*, following the SAC [18] design. The algorithms **SAC**, **DrQ**, **RAD**, and **SVEA** all use this exact architecture. The other methods include an additional component:

- **SODA** has an auxiliary network that tries to minimize the distance of augmented and non-augmented images in their feature space.
- **SGQN** comes with another head attached to the critic's large projection. This auxiliary component tries to minimize the difference between the masks it outputs and the masks generated by a saliency metric.
- **MaDi** introduces the Masker network at the start of the architecture. A small ConvNet comprising 3 layers, each with a kernel size of 3x3, stride 1, and padding 1 (with value 0). Each layer is followed by a ReLU activation, except for the output layer, which applies a Sigmoid activation function to ensure the generated mask values are between 0 and 1. See Figure 9.

*Mask Application.* In MaDi, the mask generated by the Masker is applied to the visual input by element-wise multiplication before it is passed to the actor and critic networks. Masks are generated and applied per frame, but we only need one forward pass through the Masker to get masks for all 3 frames of a stacked input observation. This is shown in the following code snippet:

```python
def apply_mask(self, obs):  # obs shape: (B,9,H,W)
    """Apply masks generated by one forward pass of the Masker"""

    # split the stacked frames
    frames = obs.chunk(3, dim=1)  # [(B,3,H,W), (B,3,H,W), (B,3,H,W)]

    # concat in batch dim
    frames_cat = torch.cat(frames, dim=0)  # (3B,3,H,W)

    # one forward pass through the Masker network
    masks_cat = self.masker(frames_cat)  # (3B,1,H,W)

    # split the batch dim back
    masks = masks_cat.chunk(3, dim=0)  # [(B,1,H,W), (B,1,H,W), (B,1,H,W)]

    # element-wise multiplication
    # uses broadcasting over the 3 RGB channels within 1 frame
    masked_frames = [m*f for m,f in zip(masks,frames)]  # [(B,3,H,W), (B,3,H,W), (B,3,H,W)]

    # concat in channel dim
    return torch.cat(masked_frames, dim=1)  # (B,9,H,W)
```
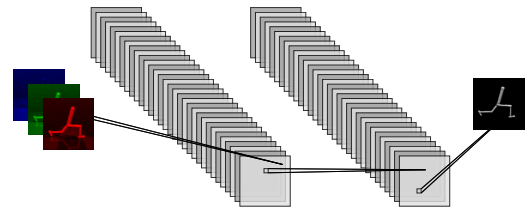


**Figure 9: The Masker network architecture contains two hidden layers with 32 channels and ReLU activations, while the output layer produces just 1 channel with a Sigmoid activation.**

## A.1 Vision Transformer implementation

The design of our Vision Transformer [8, ViT] experiments closely follows the setup of SVEA [21], where the ConvNet encoder of Figure 1 is replaced by a ViT encoder, while introducing minimal changes to the rest of the architecture. For this ViT encoder, we use RGB frames with dimensions $96 \times 96 \times 3$ as input instead of the $84 \times 84 \times 3$ used for the ConvNet. From this, we extract a total of 144 non-overlapping $8 \times 8 \times 3k$ image patches (with $k$ being the frame count in a stacked observation as described in Section 3). Each patch is mapped to a 128-dimensional embedding, which becomes a token to be processed by the ViT encoder. Consistent with the original ViT design, our model incorporates learned positional encodings and a modifiable `class` token. The architecture comprises four Transformer [45] blocks. Each block applies Multi-Head Attention encompassing eight heads, after which an MLP layer with hidden dimension 128 and a GELU [24] activation follows. We optimize the ViT encoder with the critic loss, just like the original ConvNet. We maintain the same optimizer (Adam) and learning rates as in Table 4, but we double the batch size to 256. Our ViT encoder starts without any pre-trained parameters, and we do not use weight decay. The entire training process spans around three days on a single NVIDIA RTX A4000 GPU.

## B ADDITIONAL RESULTS

This section provides additional results from our experiments, demonstrating the performance of MaDi in the original training environments and various other environments with differing levels of distractions. The table with results on `video_easy` from the DMControl Generalization Benchmark [22] are already presented in Table 1 of the main text, see Section 5.2. Here we show the results on the original training environments in Sections B.1 and B.3, as well as evaluations on `video_hard` in Section B.5, and the Distracting Control Suite in Section B.6.

*Statistical tests.* We verify the statistical significance of the results in our tables by running the unequal variances $t$-test (also called Welch's $t$-test) between MaDi and the best baseline for each environment. This test does not assume similar variances, unlike the Student's $t$-test. All our tests are two-sided, and we consider p-values less than 0.05 as statistically significant.

*Training directly on distractions.* Interestingly enough, training directly on the challenging environments with video backgrounds does not work well. In our preliminary experiments, this holds for all algorithms in both the simulation and robotic environments. Some clean data seems necessary to learn an adequate policy, before agents can learn to deal with distractions. Finding the ideal curriculum could be an interesting direction for future work.

## B.1 Results on training environment UR5 Robotic Arm

In the UR5 Robotic Arm experiments, the algorithms are trained asynchronously. This means that MaDi makes fewer updates than the others, as it takes some extra time to use the Masker network. However, even with fewer updates it surpasses the performance of other baselines on the training task; see Figure 10. The corresponding results on the generalization task are presented in Section 6, Figure 7.
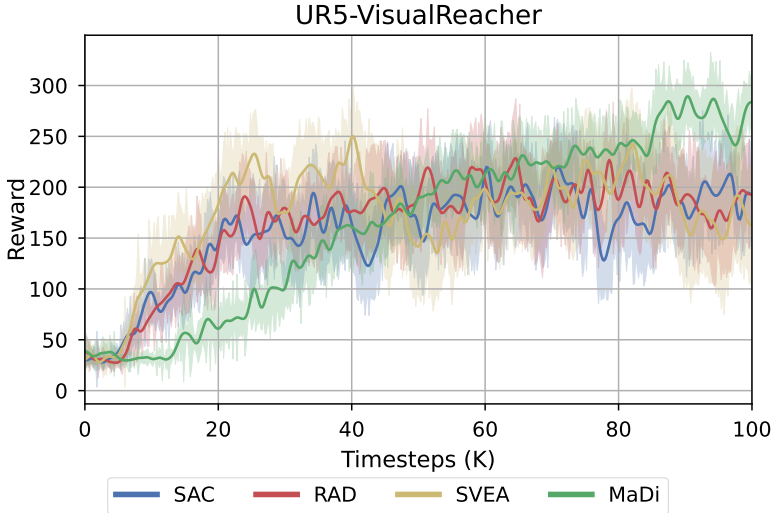


Figure 10: Performance of MaDi and three baselines on the UR5-VisualReacher training environment. Agents are trained on this clean setup (white background with the red dot) for $100K$ steps. Curves show the mean over five seeds with standard error shaded alongside. Results on the generalization task are shown in Figure 7.

## B.2 Results on UR5 Robotic Arm with *sparse* rewards

In the *sparse reward* setting of the UR5-VisualReacher task, the agent is only rewarded with a +1 whenever its camera is close enough to the red target (according to a predefined threshold). For all other situations, it receives zero reward. Even in this more challenging setting, MaDi is able to outperform the baselines in terms of generalization, as shown in Figure 11a. On the training environment, presented in Figure 11b, MaDi learns a bit slower but eventually reaches similar performance to SAC and RAD. SVEA seems to struggle a lot with sparse rewards. The MaDi and SVEA agents both use the default `overlay` augmentation here.



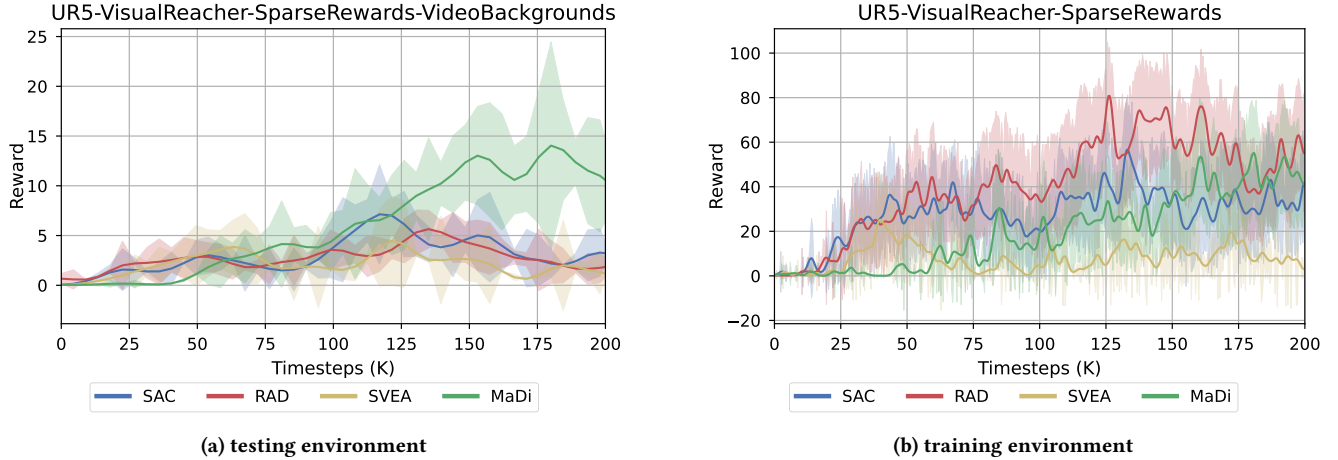(a) testing environment          (b) training environment

**Figure 11: Performance of MaDi and three baselines on the *sparse reward* version of the UR5-VisualReacher environment. Agents are trained on the clean setup (white background with the red dot) for $200K$ steps. Curves show the mean over five seeds with standard error shaded alongside.**

## B.3 Results on training environment DMControl

The results of MaDi and multiple baselines on the original training environments of the DeepMind Control Suite are shown in Table 5. The best performance is achieved by DrQ, an algorithm that only uses light augmentations (random shifts). The stronger augmentations that SODA, SVEA, SGQN, and MaDi use do not improve the results on the training environment, as they are specifically designed to enhance generalization to other domains. MaDi does still come close to the state-of-the-art in the training environments as well.

**Table 5: Performance on the original training environments of MaDi and various baseline algorithms after training for $500K$ steps. Showing mean undiscounted return and standard error over five seeds. MaDi is competitive (no significant difference) with the best baseline in 5/6 cases, denoted by ★.**

| clean | SAC | DrQ | RAD | SODA | SVEA | SGQN | MaDi | p-value |
|---|---|---|---|---|---|---|---|---|
| ball_in_cup catch | 727 ±117 | **973 ±4** | 651 ±178 | 965 ±4 | 811 ±144 | 778 ±174 | 808 ±144 | 0.313★ |
| cartpole balance | **996 ±1** | 993 ±3 | 985 ±4 | 981 ±4 | 978 ±2 | 966 ±7 | 986 ±3 | 0.024 |
| cartpole swingup | 861 ±8 | **865 ±4** | 861 ±8 | 230 ±142 | 842 ±8 | 806 ±14 | 851 ±5 | 0.077★ |
| finger spin | 324 ±9 | 643 ±36 | 593 ±13 | 493 ±117 | **711 ±55** | 597 ±13 | 681 ±17 | 0.626★ |
| walker stand | 580 ±117 | **979 ±4** | 965 ±3 | 147 ±17 | 973 ±2 | 694 ±144 | 972 ±2 | 0.147★ |
| walker walk | 366 ±42 | 892 ±33 | 837 ±21 | 496 ±173 | **937 ±13** | 892 ±27 | 893 ±27 | 0.193★ |
| avg | 642 | 891 | 816 | 552 | 875 | 789 | 865 | |

## B.4 Results on `video_easy`

The performance of MaDi and various other baseline algorithms in the `video_easy` environment are presented in Table 6. This is the same as Table 1 in the main body, but here we have enough space to include p-values of our statistical tests between MaDi and the best baseline. The `video_easy` environment represents a scenario with a low level of visual distractions, making it a suitable benchmark to test the generalization performance of reinforcement learning algorithms in a still relatively clean environment.

Table 6: Generalization performance of MaDi and various baseline algorithms on six different environments after training for $500K$ steps. Evaluated on `video_easy` from the DMControl-GB benchmark, showing mean and standard error over five seeds. MaDi significantly outperforms (★★) or is competitive to (★) the state-of-the-art in 5/6 environments.

| video_easy | SAC | DrQ | RAD | SODA | SVEA | SGQN | MaDi | p-value |
|---|---|---|---|---|---|---|---|---|
| ball_in_cup catch | 602 ±91 | 714 ±131 | 561 ±147 | 750 ±98 | 757 ±138 | 761 ±171 | **807** ±144 | 0.841★ |
| cartpole balance | 924 ±19 | 932 ±33 | 801 ±95 | 961 ±10 | 967 ±2 | 965 ±5 | **982** ±4 | 0.011★★ |
| cartpole swingup | 782 ±21 | 613 ±74 | 658 ±17 | 215 ±125 | 786 ±15 | 798 ±13 | **848** ±6 | 0.015★★ |
| finger spin | 227 ±26 | 543 ±50 | 479 ±65 | 429 ±100 | 645 ±39 | 592 ±11 | **679** ±17 | 0.461★ |
| walker stand | 507 ±113 | 954 ±10 | 961 ±1 | 147 ±17 | **977** ±3 | 672 ±153 | 967 ±3 | 0.041 |
| walker walk | 334 ±37 | 821 ±38 | 726 ±42 | 479 ±168 | **936** ±14 | 882 ±26 | 895 ±24 | 0.183★ |
| avg | 563 | 763 | 698 | 497 | 845 | 778 | 863 | |

## B.5 Results on `video_hard`

The performance of MaDi and various other baseline algorithms in the `video_hard` environment from the DMControl Generalization Benchmark [22] are presented in Table 7. The `video_hard` environment represents a scenario with a higher level of visual distractions, making it a suitable benchmark to test the generalization performance of reinforcement learning algorithms in an environment full of distractions. Similar to the results of Table 1 on `video_easy`, MaDi also shows the best generalization capability overall on `video_hard`.

Table 7: Generalization performance of MaDi and various baseline algorithms on six different environments trained for $500K$ steps. Evaluated on `video_hard` from the DMControl-GB benchmark, showing mean undiscounted return and standard error over five seeds. MaDi significantly outperforms (★★) or is competitive to (★) the state-of-the-art in 4/6 environments.

| video_hard | SAC | DrQ | RAD | SODA | SVEA | SGQN | MaDi | p-value |
|---|---|---|---|---|---|---|---|---|
| ball_in_cup catch | 176 ±38 | 180 ±51 | 261 ±58 | 194 ±39 | 327 ±59 | 687 ±155 | **758** ±135 | 0.737★ |
| cartpole balance | 314 ±12 | 250 ±6 | 263 ±18 | 449 ±46 | 579 ±26 | 703 ±17 | **827** ±25 | 0.005★★ |
| cartpole swingup | 140 ±10 | 155 ±10 | 133 ±19 | 172 ±46 | 453 ±26 | 488 ±18 | **619** ±24 | 0.003★★ |
| finger spin | 21 ±4 | 16 ±3 | 8 ±2 | 107 ±30 | 154 ±31 | **554** ±8 | 358 ±25 | 0.001 |
| walker stand | 233 ±28 | 361 ±44 | 343 ±31 | 143 ±12 | 847 ±18 | 606 ±92 | **920** ±14 | 0.012★★ |
| walker walk | 168 ±19 | 111 ±16 | 116 ±18 | 214 ±68 | 526 ±55 | **718** ±37 | 504 ±33 | 0.003 |
| avg | 175 | 179 | 187 | 213 | 481 | 626 | 664 | |

## B.6 Results on DistractingCS

We demonstrate the performance of MaDi and our baseline algorithms in the Distracting Control Suite [39] with intensity 0.1 in Table 8. DistractingCS is an extreme benchmark that features different distractions, such as varying camera angles, changing agent colors, and random videos in the background. The results indicate that MaDi performs well on this very challenging benchmark.

Table 8: Generalization performance of MaDi and various baseline algorithms on six different environments after training for $500K$ steps. We evaluate on the Distracting Control Suite with intensity $0.1$, showing mean undiscounted return and standard error over five seeds. MaDi is competitive (no significant difference) with the best baseline in 5/6 environments, denoted by ⋆.

| distracting_cs | SAC | DrQ | RAD | SODA | SVEA | SGQN | MaDi | p-value |
|---|---|---|---|---|---|---|---|---|
| ball_in_cup catch | 112 ±69 | 177 ±67 | 33 ±30 | 107 ±30 | 105 ±44 | **529** **±111** | 477 ±99 | 0.738⋆ |
| cartpole balance | 310 ±22 | 237 ±19 | 233 ±19 | 210 ±6 | 278 ±22 | 358 ±17 | **377** **±12** | 0.373⋆ |
| cartpole swingup | 146 ±8 | 158 ±13 | 162 ±6 | 138 ±8 | 219 ±18 | **283** **±20** | 252 ±30 | 0.412⋆ |
| finger spin | 28 ±5 | 55 ±13 | 38 ±15 | 52 ±18 | 137 ±21 | **386** **±13** | 211 ±30 | 0.002 |
| walker stand | 353 ±66 | 727 ±58 | 539 ±43 | 157 ±20 | **871** **±22** | 641 ±110 | 866 ±16 | 0.875⋆ |
| walker walk | 220 ±12 | 429 ±19 | 291 ±30 | 235 ±91 | 568 ±28 | **642** **±44** | 570 ±49 | 0.303⋆ |
| avg | 195 | 297 | 216 | 150 | 363 | 473 | 459 | |

## C    MASK EXAMPLES

In Figure 12, we show examples of masked observations by the Masker network of MaDi. These masked observations are seen by the actor and critic networks, enabling them to perform well even in the presence of visual distractions.
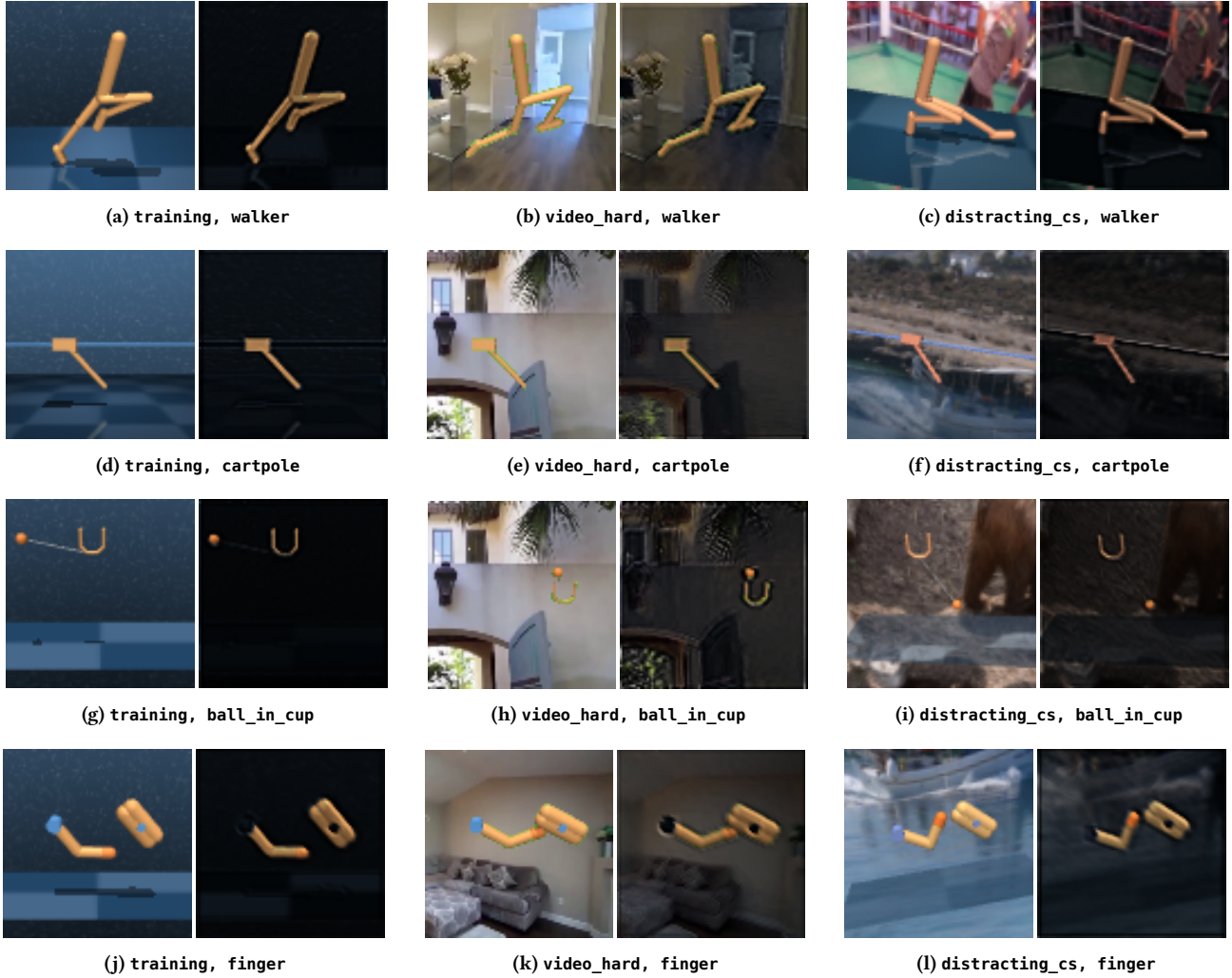


(a) `training, walker`

(b) `video_hard, walker`

(c) `distracting_cs, walker`

(d) `training, cartpole`

(e) `video_hard, cartpole`

(f) `distracting_cs, cartpole`

(g) `training, ball_in_cup`

(h) `video_hard, ball_in_cup`

(i) `distracting_cs, ball_in_cup`

(j) `training, finger`

(k) `video_hard, finger`

(l) `distracting_cs, finger`

Figure 12: Examples of masked observations by MaDi under the `overlay` augmentation in training and testing environments (`video_hard` and `distracting_cs`) for all domains used. The Masker network produces useful masks, displaying task-relevant information and dimming the distractions.

## C.1 Masks over time

In this section we show how the masks evolve over time. The Masker network is initialized in such a way that the first masks consist of values very close to 0.5, resulting in a fully gray image. Then, after just a small number of updates, the Masker can already generate useful masks that focus on the relevant features. See Figures 13 and 14.
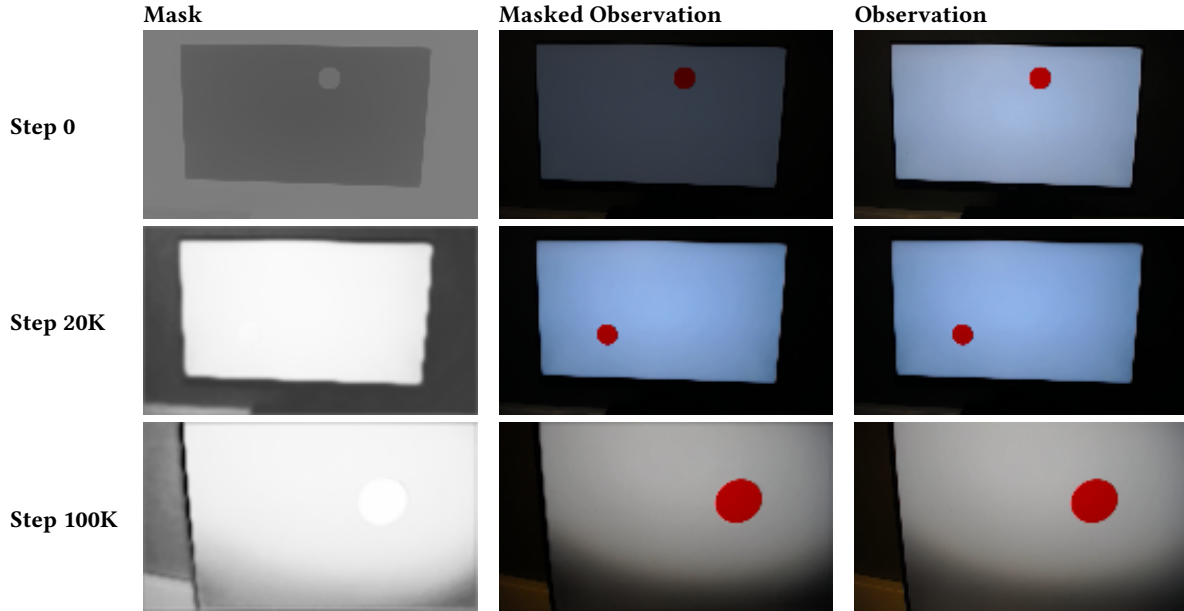


Figure 13: Masks produced by MaDi (with `overlay` augmentation) on the UR-VisualReacher training environment. After just 20,000 timesteps MaDi has learned to output a subtle, but useful mask. The initial Masker outputs values close to 0.5, shown by the gray image in the top-left corner (vaguely revealing the shape of the target circle).



Figure 14: Masks produced by MaDi (with `overlay` augmentation) on the UR-VisualReacher-VideoBackgrounds testing environment. After 20,000 timesteps MaDi has learned to output a subtle, but helpful mask.

# D   OTHER AUGMENTATIONS

In this section, we present the results of MaDi under different augmentation techniques instead of the default `overlay`. We experimented with a random convolutional layer as augmentation, called `conv` in [22], and an augmentation based on hue-saturation-value (HSV) thresholds named `splice` [55]. In Figure 15, we present an overview of the augmentations used.

MaDi can perform quite well with the `splice` augmentation. The masks often look even clearer than with `overlay`, as shown in Figure 16. However, we have also occasionally seen masks that turn fully-black with the `splice` augmentation, which hides all relevant information needed by the actor to take an optimal action. This is likely why the performance of MaDi with `splice` is not the best in all environments. Further research is necessary to determine why MaDi sometimes produces black masks under this augmentation technique.
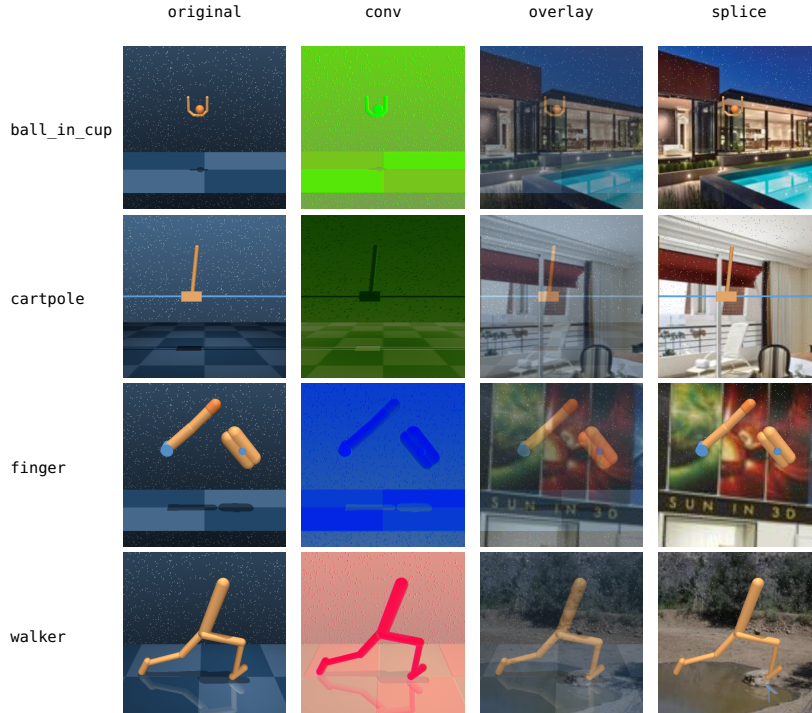


**Figure 15: Various augmentations on the four domains of the DeepMind Control Suite used in our study. Applying `overlay` results in dark-blue and blurry visuals, whereas `splice` produces crisper images. The `splice` augmentation requires some prior knowledge about the environment to set the HSV thresholds.**

**Table 9: Comparing the effect of different data augmentations. We show the generalization performance of MaDi on multiple benchmarks, using three different augmentations. The agents are all trained for $500K$ steps. We present the mean undiscounted return and standard error over five seeds.**

### (a) video_easy

| augmentation | overlay | conv | splice |
|---|---|---|---|
| ball_in_cup catch | **807** ±144 | 603 ±115 | 617 ±196 |
| cartpole balance | **982** ±4 | 903 ±37 | 881 ±91 |
| cartpole swingup | **848** ±6 | 693 ±34 | 785 ±22 |
| finger spin | **679** ±17 | 304 ±85 | 568 ±131 |
| walker stand | **967** ±3 | 408 ±149 | 828 ±129 |
| walker walk | **895** ±24 | 476 ±60 | 861 ±10 |
| avg | **863** | 565 | 757 |

### (b) video_hard

| augmentation | overlay | conv | splice |
|---|---|---|---|
| ball_in_cup catch | **758** ±135 | 84 ±30 | 610 ±194 |
| cartpole balance | 827 ±25 | 320 ±34 | **879** ±90 |
| cartpole swingup | 619 ±24 | 161 ±14 | **769** ±27 |
| finger spin | 358 ±25 | 61 ±18 | **566** ±130 |
| walker stand | **920** ±14 | 210 ±29 | 824 ±128 |
| walker walk | 504 ±33 | 99 ±12 | **858** ±17 |
| avg | 664 | 156 | **751** |

### (c) distracting_cs

| augmentation | overlay | conv | splice |
|---|---|---|---|
| ball_in_cup catch | 477 ±99 | 165 ±61 | **518** ±190 |
| cartpole balance | 377 ±12 | 241 ±8 | **409** ±24 |
| cartpole swingup | 252 ±30 | 176 ±8 | **300** ±23 |
| finger spin | **211** ±30 | 111 ±42 | 72 ±20 |
| walker stand | **866** ±16 | 271 ±56 | 346 ±59 |
| walker walk | **570** ±49 | 378 ±31 | 230 ±28 |
| avg | **459** | 224 | 313 |

(a) `training, walker`      (b) `video_hard, walker`      (c) `distracting_cs, walker`

(d) `training, cartpole`      (e) `video_hard, cartpole`      (f) `distracting_cs, cartpole`

(g) `training, ball_in_cup`      (h) `video_hard, ball_in_cup`      (i) `distracting_cs, ball_in_cup`
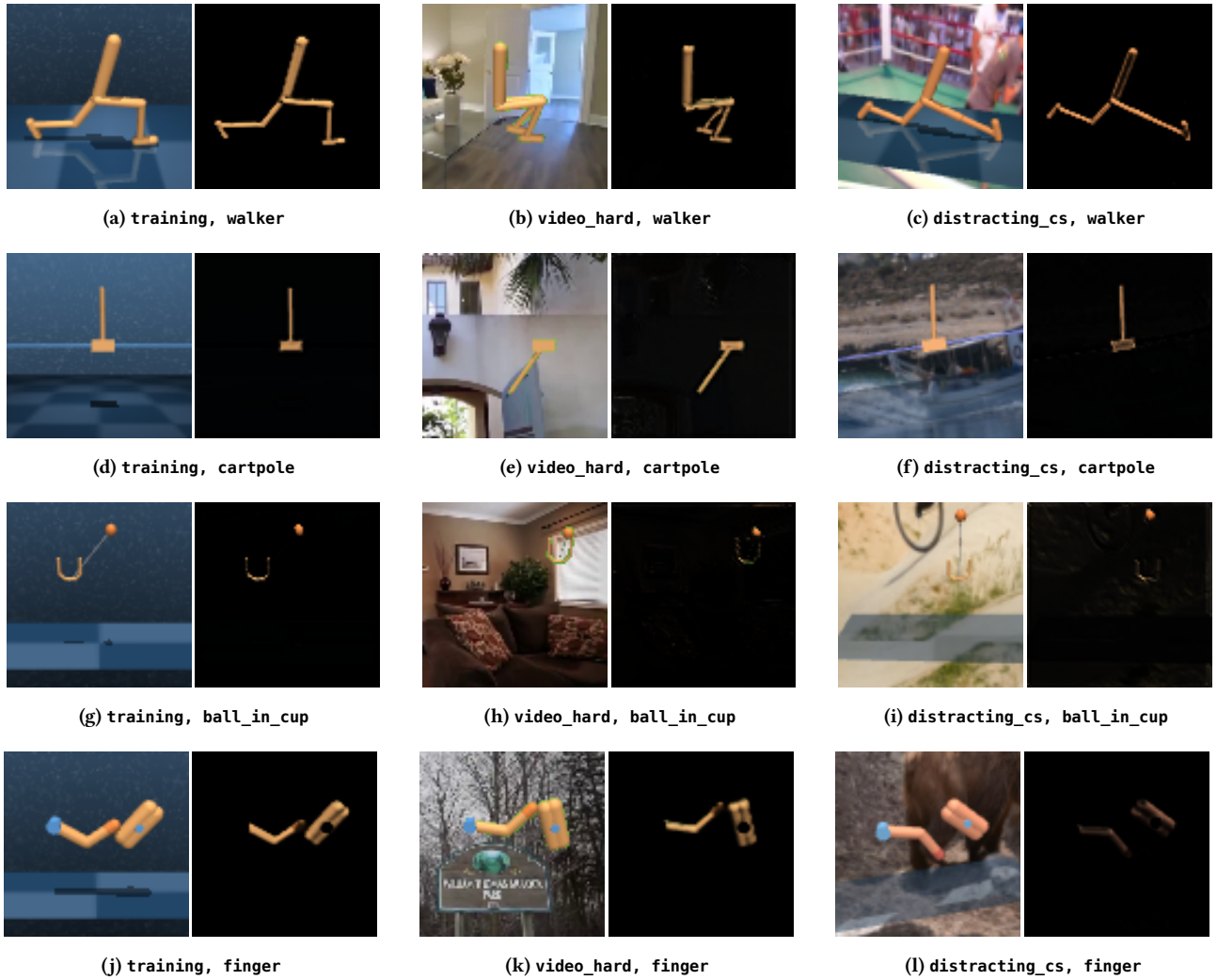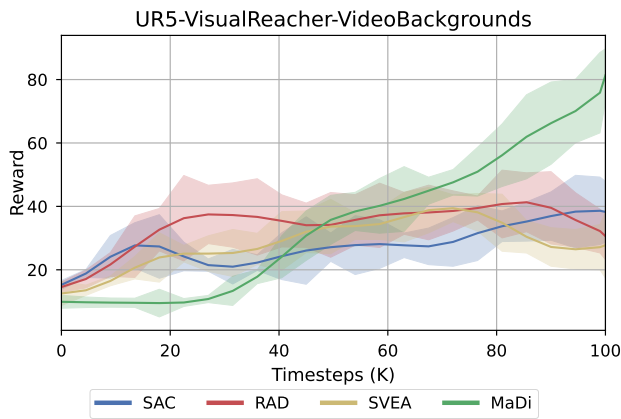
(j) `training, finger`      (k) `video_hard, finger`      (l) `distracting_cs, finger`

**Figure 16: Examples of masked observations by MaDi under the `splice` augmentation in training and testing environments (`video_hard` and `distracting_cs`) for all domains used. The Masker network generally produces clear and precise masks, but it can sometimes mask too much (i.e., parts of the agent) when using `splice`, as seen in the bottom right.**
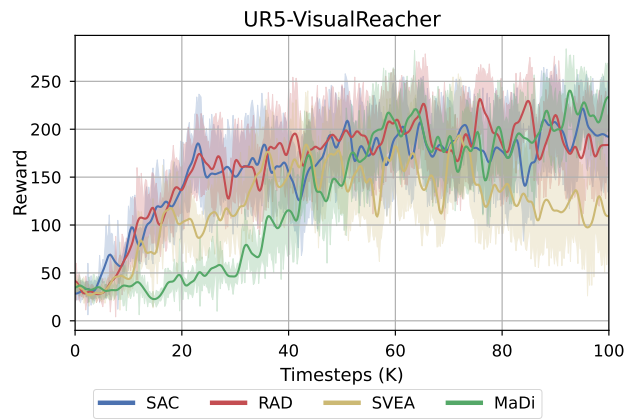
## D.1 Robotic experiments with conv augmentation

We have also run experiments with the `conv` augmentation on our UR-VisualReacher task, to see whether it can help to improve generalization on the video backgrounds, see Figure 17. The experiments presented in the main body use the `overlay` augmentation. (These results are shown in Figure 7 and the corresponding training environment curves in Figure 10.)

The masks generated by MaDi under the `conv` augmentation seem to be a little clearer than the masks generated under the `overlay` augmentation. We think that MaDi with `overlay` was still able to perform better as the encoder network (see Figure 1) also benefits from seeing some strong `overlay` augmentation during training. Especially early on, when the masks are not fully formed yet, the encoder will receive quite a lot of augmented pixels, which can help its robustness. In Figures 18 and 19 we show some examples of masks that MaDi creates under the `conv` augmentation. These are all from one particular seed, masks from other seeds are generally similar but can differ slightly in their clarity.

**(a) testing environment**

**(b) training environment**

Figure 17: Performance of MaDi and three baselines on the UR5-VisualReacher environment, where both MaDi and SVEA use the **conv** augmentation now. Agents are trained on the clean setup (white background with the red dot) for $100K$ steps. Curves show the mean over five seeds with standard error shaded alongside.
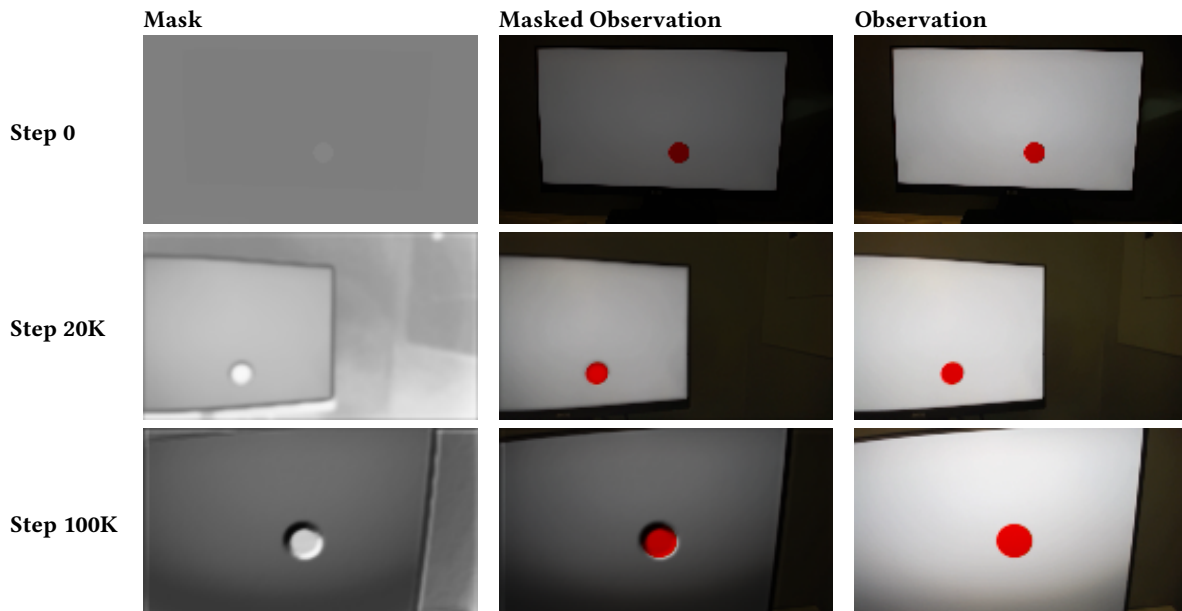


Figure 18: Masks produced by MaDi (with **conv** augmentation) on the UR-VisualReacher training environment. After just $20,000$ timesteps MaDi has already learned to output useful masks. The initial Masker outputs values close to $0.5$, shown by the gray image in the top-left corner (which vaguely reveals the shape of the target circle).
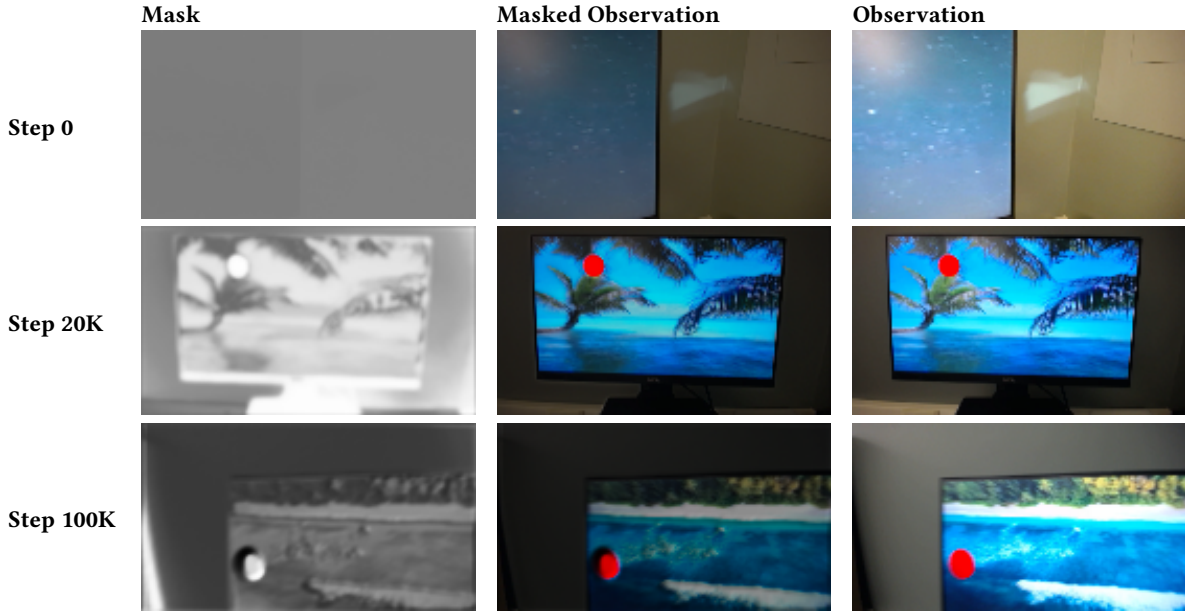
**Figure 19: Masks produced by MaDi (with conv augmentation) on the UR-VisualReacher-VideoBackgrounds testing environment. After just** $20,000$ **timesteps MaDi can generate useful masks on the unseen testing environment. The initial Masker network outputs values close to** $0.5$, **shown by the gray image in the top-left corner.**

## E HARDWARE RESOURCES

*Simulation experiments.* The experiments of this study were all individually run on single GPU setups. We have operated on distinct Compute Clusters with as GPU types: NVIDIA RTX A4000, NVIDIA V100SXM2, and NVIDIA P100 Pascal, with access to $4 - 6$ CPU cores each. A full training run of 500K timesteps can take up to 36 hours with batch size 128. Most of our exploratory runs were done with batch size 64. This halves the runtime approximately, while begetting little degradation in performance.

*Robotic experiments.* We use a workstation with an AMD Ryzen Threadripper 2950 processor, an NVIDIA 2080Ti GPU, and 128G memory for the UR5 Robotic Arm experiments. The software setup uses multi-processing and multi-threading extensively to communicate between the arm, the camera, and the workstation. A full training run of 100K timesteps takes around 2.5 hours. In the sparse reward setting we train for 200K steps, requiring 5 hours. The batch size has no impact here, as we update the neural networks asynchronously in this setup.

## F TASK DESCRIPTIONS

In this section we provide additional details on the particular environments used in this study, both in simulation and on a physical robot.

### F.1 Simulation experiments

We experiment on tasks from the benchmarks DMControl-GB [22] and DistractingCS [39], both of which build on environments from DMControl [41]. We clarify and provide properties of all the tasks considered in our study below, based on the descriptions from [21, 41]. These DMControl tasks are selected based on previous work on both sample efficiency [19, 50] and generalization [22, 39], and represent a diverse and challenging skill set in the context of image-based RL. All environments emit observation frames $\mathbf{o} \in \mathbb{R}^{84 \times 84 \times 3}$. Three frames are stacked together to make one state $\mathbf{s} \in \mathbb{R}^{84 \times 84 \times 9}$. See Figure 12 for examples of observations from each domain.

- `ball_in_cup-catch` ($\mathbf{a} \in \mathbb{R}^2$). An actuated planar receptacle is to move around in the plane and catch a ball attached by a string to its bottom. Sparse rewards.
- `cartpole-balance` ($\mathbf{a} \in \mathbb{R}$). Balance an unactuated pole by moving a cart left or right at the base of the pole. The pole starts in upright position. The agent is rewarded for balancing the pole within a fixed angle. Dense rewards.
- `cartpole-swingup` ($\mathbf{a} \in \mathbb{R}$). Swing up and balance an unactuated pole by moving a cart left or right at the base of the pole. The pole starts in downwards orientation. The agent is rewarded for balancing the pole within a fixed angle. Dense rewards.

- `finger-spin` ($a \in \mathbb{R}^2$). A manipulation problem with a planar 2 DoF finger. The task is to continually rotate a free body on an unactuated hinge. Sparse rewards.
- `walker-stand` ($a \in \mathbb{R}^6$). A bipedal planar walker that is rewarded for standing with an upright torso at a constant minimum height. Dense rewards.
- `walker-walk` ($a \in \mathbb{R}^6$). A bipedal planar walker that is rewarded for walking forward at a certain target velocity. Dense rewards.

## F.2 Robotic experiments

We train MaDi and three other baselines on the UR-VisualReacher task [48, 52] shown in Figure 20, and test the generalization performance on a novel experiment that we designed, dubbed UR-VisualReacher-VideoBackgrounds. The task of the robotic arm is the same in both environments: it gets rewarded for positioning its camera as close as possible to the red dot on the screen. In the testing environment, one of the videos shown in Figure 21 is playing in the background. We take videos from the DeepMind Control Generalization Benchmark [22].

For the reward function: we adjust the RGB thresholds used to determine whether pixels are part of the red dot or not, such that it is robust against the video backgrounds. We selected videos that do not have any strong red colors in them, so that the rewards are well calibrated. It would be an interesting direction for future work to verify whether MaDi can detect the red dot even if other task-irrelevant pixels are red too. This would require a different setup of the reward function.

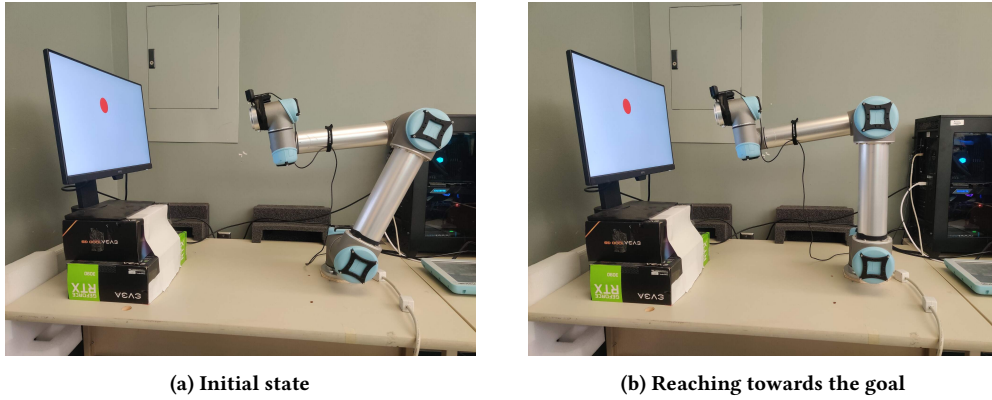

(a) Initial state    (b) Reaching towards the goal

Figure 20: The UR5 Robotic Arm setup used in our experiments. The agent is rewarded for getting its camera as close as possible to the red dot randomly located on the screen.
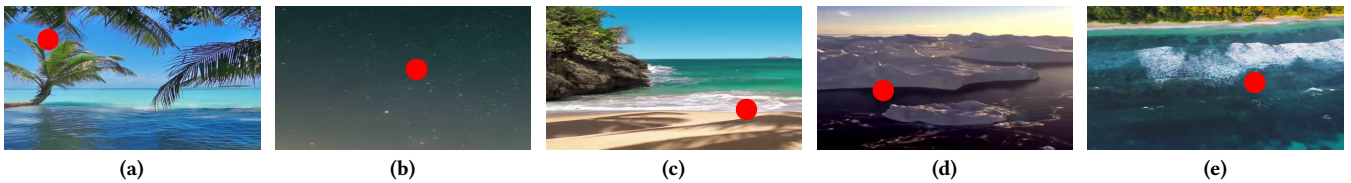


(a)    (b)    (c)    (d)    (e)

Figure 21: Examples of frames from each video shown on the screen in the UR5-VisualReacher task to test generalization.