# Simple Sensor Intentions for Exploration

Tim Hertweck*, Martin Riedmiller*, Michael Bloesch, Jost Tobias Springenberg,
Noah Siegel, Markus Wulfmeier, Roland Hafner, Nicolas Heess

DeepMind

London, United Kingdom

*Abstract*—Modern reinforcement learning algorithms can learn solutions to increasingly difficult control problems while at the same time reduce the amount of prior knowledge needed for their application. One of the remaining challenges is the definition of reward schemes that appropriately facilitate exploration without biasing the solution in undesirable ways, and that can be implemented on real robotic systems without expensive instrumentation. In this paper we focus on a setting in which goal tasks are defined via simple sparse rewards, and exploration is facilitated via agent-internal auxiliary tasks. We introduce the idea of simple sensor intentions (SSIs) as a generic way to define auxiliary tasks. SSIs reduce the amount of prior knowledge that is required to define suitable rewards. They can further be computed directly from raw sensor streams and thus do not require expensive and possibly brittle state estimation on real systems. We demonstrate that a learning system based on these rewards can solve complex robotic tasks in simulation and in real world settings. In particular, we show that a real robotic arm can learn to grasp and lift and solve a Ball-in-a-Cup task from scratch, when only raw sensor streams are used for both controller input and in the auxiliary reward definition. A video showing the results can be found at https://deepmind.com/research/publications/Simple-Sensor-Intentions-for-Exploration.
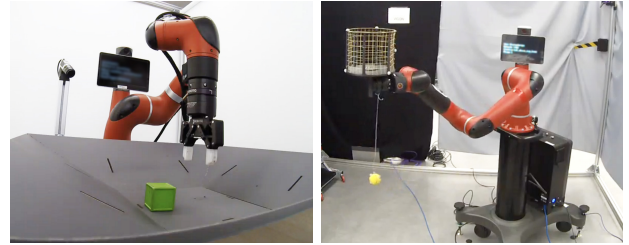
Fig. 1. Manipulation setup with a Rethink Sawyer robotic arm and a Robotiq 2F-85 parallel gripper (left). Ball-in-a-cup task setup with a Rethink Sawyer robotic arm and a custom made Ball-and-Cup attachment (right).

## I. INTRODUCTION

An important consideration on the path towards general AI is to minimize the amount of prior knowledge needed to set up learning systems. Ideally, we would like to identify principles that transfer to a wide variety of different problems without the need for manual tuning and problem-specific adjustments. In recent years, we have witnessed substantial progress in the ability of reinforcement learning algorithms to solve difficult control problems from first principles, often directly from raw sensor signals such as camera images, without need for carefully handcrafted features that would require human understanding of the particular environment or task [27, 40].

One of the remaining challenges is the definition of reward schemes that appropriately indicate task success, facilitate exploration without biasing the solution in undesirable ways, and that can be implemented on real robotics systems without expensive instrumentation. In this paper we focus on situations, where the external task is given by a sparse reward signal, that is '1' if and only if the task is solved. Such reward functions are often easy to define, and by solely focusing on task success strongly mitigate the bias on the final solution. The associated challenge is, however, that starting from scratch with a naive exploration strategy will most likely never lead

to task success and the agent will hardly receive any learning signal.

The Scheduled Auxiliary Control (SAC-X) [33] framework tackles this problem by introducing a set of auxiliary rewards, that help the agent to explore the environment. For each auxiliary reward an auxiliary policy ('intention') is learned and executed to collect data into a shared replay-buffer. This diverse data facilitates learning of the main task. The important insight of [33] is that the exact definition of auxiliary tasks can vary, as long as they jointly lead to an adequate exploration strategy that allows to collect rich enough data such that learning of the main task can proceed. However, in the original work the auxiliary tasks are still defined with some semantic understanding of the environment in mind, e.g. 'move an object' or 'place objects close to each other'. This requires both task-specific knowledge for the definition of the auxiliaries, as well as the technical prerequisites to estimate the relevant quantities required for the computation of the rewards - e.g. camera calibration, object detection and object pose estimation. In this work we want to make a step towards a more generic approach for defining auxiliary tasks, that reduces the need for task-specific semantic interpretation of environment sensors, in particular of camera images.

A fundamental principle to enable exploration in sparse reward scenarios, is to learn auxiliary behaviours that deliberately change sensor responses. While variants of this idea have been already suggested earlier, e.g. [42, 23], we here introduce a generic way to implement this concept into the SAC-X framework: 'simple sensor intentions' (SSIs) encourage the agent to explore the environment and help with collecting meaningful data for solving the sparsely rewarded main task. Being largely task agnostic, SSIs can be reused across tasks with no or only minimal adjustments. Further, simple sensor

*Equal contribution. Correspondence to: thertweck@google.com, riedmiller@google.com.

intentions (SSIs) are based on rewards defined on the deliberate change of scalar sensor responses, that are derived from raw sensor values. We propose two ways to effect change, namely (a) to attain certain set points (like e.g. the minimum or maximum of a sensor response), or (b) by rewarding increase or decrease of a sensor response.

However, not all sensory inputs can be directly mapped to scalar values, like e.g. raw camera images. As an exemplary procedure, we suggest a concrete pre-processing for mapping raw images into simple sensor responses by computing and evaluating basic statistics, such as the spatial mean of color-filtered images. While SSIs propose a general way to deal with all kinds of sensor values available in a robotic system (like touch sensors, joint angle sensors, position sensors, ...), we mostly investigate pixel inputs here as an example of a sensor type that is widely used in robotics.

As a proof of concept, we show that these simple sensor intentions (SSIs) can be applied on a variety of interesting robotic domains, both in simulation and on real robots. Most notably, we show that with SSIs, SAC-X is capable of learning to play the Ball-in-a-Cup game from scratch, purely from pixel and proprioceptive inputs - for both as observation and for computing the auxiliary rewards given to the agent.

## II. PRELIMINARIES

We consider a reinforcement learning setting with an agent operating in a Markov Decision Process (MDP) consisting of the state space $\mathcal{S}$, the action space $\mathcal{A}$ and the transition probability $p(s_{t+1}|s_t, a_t)$ of reaching state $s_{t+1}$ from state $s_t$ when executing action $a_t$ at the previous time step $t$. The goal of an agent in this setting is to succeed at learning a given task $k$ out of a set of possible tasks $\mathcal{K}$. The actions are assumed to be drawn from a probability distribution over actions $\pi_k(a|s)$ referred to as the agent's policy for task $k$. After executing an action $a_t$ in state $s_t$ and reaching state $s_{t+1}$ the agent receives a task-dependent, scalar reward $r_k(s_t, s_{t+1})$. Given a target task $g$ we define the expected return (or value) when following the task-conditioned policy $\pi_g$, starting from state $s$, as

$$V^{\pi_g}(s) = \mathbb{E}_{\pi_g}\big[\sum_{t=0}^{\infty} \gamma^t r_g(s_t, s_{t+1}) \mid s_0 = s\,\big],$$

with $a_t \sim \pi_g(\cdot|s_t)$ and $s_{t+1} \sim p(\cdot|s_t, a_t)$ for all $s \in \mathcal{S}$. The goal of Reinforcement Learning then is to find the policy $\pi_g^*$ that maximizes the value. The auxiliary intentions, defined in the following based on simple sensor rewards – that is task $k \in \mathcal{K}, k \neq g$ – give rise to their own values and policies and serve as means for efficiently exploring the MDP.

## III. SIMPLE SENSOR INTENTIONS

The key idea behind simple sensor intentions can be summarized by the following principle: in the absence of an external reward signal, a sensible exploration strategy can be formed by learning policies that deliberately cause an effect on the observed sensor values by e.g. driving the sensor responses to their extremas or by controling the sensor readings at desirable set-points.

Clearly, policies that achieve the above can learn to cover large parts of the observable state-space, even without external rewards, and should thus be useful for finding 'interesting' regions in the state-space. This idea is distinct from, but reminiscent of, existing work on intrinsic motivation for reinforcement learning [18] which often defines some form of curiosity (or coverage) signal that is added to the external reward during RL. Our goal here is to learn separate exploration policies from general auxiliary tasks, that can collect good data for the main learning task at hand.

As a motivating example, an analogy to this idea can be found by considering the exploration process of an infant: in the absence of 'knowing what the world is about', a baby will often move its body in a seemingly 'pointless', but goal directed, way until it detects a new stimulus in its sensors, e.g. a new sense of touch at the fingertips or a detected movement in a toy.

Simple sensor intentions (SSIs) propose a generic way to implement this principle in the multi-task agent framework SAC-X. SSIs are a set of auxiliary tasks, defined by standardized rewards over a set of scalar sensor responses. While many sensors in robotics naturally fit into this scheme (like e.g. a binary touch sensor or a sensor for a joint position), other sensors like raw camera images may need some transformation first in order to provide a scalar sensor signal that can be used to define a reasonable simple sensor intention.

In general, SSIs are derived from raw sensor observations in two steps:

*First step:* In the first step we map the available observations to scalar sensor responses that we want to control. Each observation $o \in s$ is a vector of sensor values coming from different sensory sources like e.g. proprioceptive sensors, haptic sensors or raw images.

We define a scalar (virtual) sensor response by mapping an observation $o$ to a scalar value $z$, i.e.

$$z = f(o), \text{ where } o \in s,$$

and where $f$ is a simple transformation of the observation. For scalar sensors, $f$ can be the identity function, while other sensors might require some pre-processing (like e.g. raw camera images, for which we describe a simple transformation in more detail in section IV). In addition – as a consequence of choosing the transformation – for each derived scalar sensor value $z$ we can also assume to know the maximum and minimum attainable value $[z_{\min}, z_{\max}]$.

*Second step:* In the second step we calculate a reward following one of two different schemes (described in detail below):

  a) rewarding the agent for reaching a specific target response, or
  b) rewarding the agent for incurring a specific change in response.

Importantly, these schemes do not require a detailed semantic understanding of the environment: changes in the environment

may have an a-priori unknown effect on sensor responses. Irregardless of this relationship between environment and sensor values we follow the hypothesis outlined earlier: a change in a sensor response indicates some change in the environment and by learning a policy that deliberately triggers this change (a simple sensor intention, SSI) we obtain a natural way of encouraging diverse exploration of the environment.

### A. Target Response Reward

Let $z_t \in \{z_t \in \mathbb{R} \mid z_{\min} \le z_t \le z_{\max}\}$ be a sensor response $z_t = f(o)$ for observation $o \in s_t$. We define the reward at time step $t$ for controlling the response $z$ towards a desired set point $\hat{z}$ as

$$r^{\hat{z}}(s_t) := 1 - \frac{|z_t - \hat{z}|}{z_{\max} - z_{\min}},$$

where $\hat{z}$ is the set point to be reached. The set point could be chosen arbitrarily in the range $[z_{\min}, z_{\max}]$, a sensible choice that we employ in our experiments, is to use the minimum and maximum response values as set points to encourage coverage of the sensor response space. We denote the two corresponding rewards as *'minimize z'* and *'maximize z'*, respectively in our experiments.

### B. Response Change Reward

While set point rewards encourage controlling sensor values to a specific value, response change rewards encourage the policy to incur a signed change of the response. Let $\Delta_t^z := z_t - z_{t-1}$ be the temporal difference between consecutive responses. We define the reward for *changing* a sensor value $z$ as

$$r^{\Delta}(s_t) := \frac{\alpha \Delta_t^z}{z_{\max} - z_{\min}}.$$

where $\alpha \in \{1, -1\}$ serves to distinguish between *increasing* and *decreasing* sensor responses. In both cases, undesired changes are penalized by our definition. This ensures that a successful policy moves the response consistently in the desired direction, instead of exploiting positive rewards by developing a cyclic behaviour [32, 29].

## IV. Transforming Images to Simple Sensor Intentions

In the following, we describe the transformation we use to obtain scalar sensor responses from raw camera images. Cameras are an especially versatile and rich type of sensor for observing an agent's surroundings and are thus particularly interesting for our purposes. Cameras typically deliver two dimensional pixel intensity arrays as sensor values. There are numerous ways to map these pixel arrays to one dimensional sensor responses, which can then be used as part of the simple sensor reward schemes.

In principle, one could treat every pixel channel as an individual response, or calculate averages in regions of the image (similar to 'pixel control' in [23]), and subsequently define sensor rewards for each of the regions. However, our
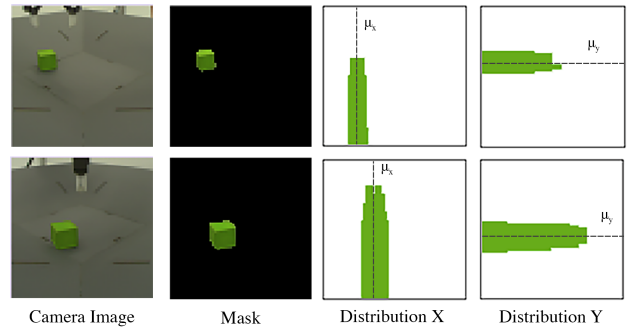


Fig. 2. The transformation used for deriving one dimensional sensor responses from camera images. We compute a binary mask and it's spatial distribution along the image axes and select the resulting distribution's mean as a sensor response.

goal is to learn a sensor intention policy for e.g. maximizing / minimizing each sensor response, that can then be executed to collect data for a target task. In such a setting, having a smaller amount of efficient exploration policies is preferable (opposed to the large number of policies a reward based on single pixel values would mandate). We therefore propose to transform images into a small amount of sensor responses by aggregating statistics of an image's spatial color distribution. As illustrated in Figure 2, we first threshold the image to retain only a given color (resulting in a binary mask) and then calculate the mean location of the mask along each of the image's axes, which we use as the sensor value. Formally we can define the two corresponding sensor values for each camera image as

$$z_x^{c_{\text{range}}} = \frac{1}{W} \sum_{x=0}^{W} x \max_y [\mathbf{1}_{c_{\text{range}}}(o_{\text{image}})[y,x]]$$

$$z_y^{c_{\text{range}}} = \frac{1}{H} \sum_{x=0}^{H} y \max_x [\mathbf{1}_{c_{\text{range}}}(o_{\text{image}})[y,x]],$$

where $H$ denotes the image height and $W$ the width, and $c_{\text{range}} = [c_{\min}, c_{\max}]$ correspond to color ranges that should be filtered. Combined with the reward scheme above, these simple sensor responses can result in intentions that try to color given positions of an image in a given color. Perhaps surprisingly, we find such a simple reward to be sufficient for encouraging meaningful exploration in a range of tasks in the experiments. We note that the reward formulation outlined in Section III mandates a sensor response to be available at each time step $t$. To avoid issues in case no pixel in the image matches the defined color range, we set any reward based on $z^{\text{image}}$ to zero at time $t = 0$ and subsequently always fall back to the last known value for the reward if no pixel response matches.

The choice of the 'right' color range $c_{\text{range}}$, for a task of interest, is a design decision that needs to be made manually. In practice, we define a set of color ranges (and corresponding sensor values) from rough estimates of the color of objects of interest in the scene. Alternatively, the color filters could be defined very broadly, which increases generality and transfer. For example, similar to a baby's preference for bright or vivid colors, one might use a color filter that matches a

broad range of hues but only in a narrow saturation range. Furthermore, if the number of interesting color ranges is large (or one chooses them randomly) one can also define aggregate intentions, where sensor rewards are averaged over several sensor responses for various color channels, such that the resulting intention has the goal of changing an arbitrary color channel's mean instead of a specific one. In addition to manually defining color ranges, we as well conduct experiments with this aggregate approach.

## V. LEARNING SIMPLE SENSOR INTENTIONS FOR ACTIVE EXPLORATION

In general, exploration policies based on the sensor rewards described in the previous section could be learned with any multi-task Reinforcement Learning algorithm or, alternatively, added as exploration bonus for an Reinforcement Learning algorithm that optimizes the expected reward for the target task $g$.

In this work we demonstate how simple sensor intentions can be used in the context of a multi-task RL algorithm, to facilitate exploration. We make use of ideas from the recent literature on data-efficient multi-task RL with auxiliary tasks (in our case defined via the simple sensory rewards). Concretely, we follow the setup from Scheduled Auxiliary Control (SAC-X) [33] and define the following policy optimization problem over $\mathcal{K}$ tasks:

$$\arg\max_{\pi} \mathbb{E}_{s \in \mathcal{B}}\Big[\sum_{k=1}^{K} \mathbb{E}_{a \sim \pi_k(\cdot|s)}[Q_\phi^k(s,a)]\Big],$$

where $\pi_k(a|s)$ is a task-conditioned policy and $Q_\phi(s,a,k)$ is a task-conditional Q-function (with parameters $\phi$); that is learned alongside the policy by minimizing the squared temporal difference error:

$$\min_{\phi} \mathbb{E}_{(s,a,s') \in \mathcal{B}}\Big[\sum_{k=1}^{K}\big(r_k(s,a)+\gamma\mathbb{E}_{\pi_k}[Q_{\hat\phi}^k(s',a')]-Q_\phi^k(s,a)\big)^2\Big],$$

where $\hat\phi$ are the periodically updated parameters of a target network. We refer to the appendix for a detail description of the neural networks used to represent $\pi_k(a|s)$ and $Q_\phi(s,a,k)$.

The set of tasks $\mathcal{K}$ is given by the reward functions for each of the SSIs that we want to learn, as well as the externally defined goal reward $r_g$. The transition distribution, for which the policy and Q-function are learned, is obtained from the replay buffer $\mathcal{B}$, which is filled by *executing both the policy for the target task $\pi_g$ as well as all other available exploration SSIs $\pi_k \in \mathcal{K}, k \neq g$.* In SAC-X, each episode is divided into multiple sequences and a policy to execute is chosen for each of the sequences. The decision which policy to execute is either made at random (referred to as Scheduled Auxiliary Control with uniform sampling or SAC-U) or based on a learned scheduler that maximizes the likelihood of observing the sparse task reward $r_g$ (referred to as SAC-Q). More details on the Reinforcement Learning procedure can be found in the Appendix as well as in the original Scheduled Auxiliary Control publication [33].

## VI. EXPERIMENTS

In the following, simple sensor intentions (SSIs) based on basic sensors (like e.g. touch) and more complex sensors, like raw camera images, are applied to several robotic experiments in simulation and on a real robot. We show, that by using the concept of SSIs, several complex manipulation tasks can be solved: grasping and lifting an object, stacking two objects and solving a ball-in-cup task end-to-end from raw pixels. In all experiments, we assume the final task reward to be given in form of a sparse (i.e. binary) external reward signal.

### A. Experimental Setup

In all following experiments we employ a Rethink Sawyer robotic arm, with either a Robotiq 2F-85 parallel gripper as end-effector or, in the case of the Ball-in-a-Cup task, with a custom made cup attachment. In the manipulation setups, the robot faces a 20cm x 20cm basket, containing a single colored block, or - in the case of the stack task - two differently colored blocks. The basket is equipped with three cameras, that are used as the only exteroceptive inputs in all manipulation experiments (Figure 1, left). For the Ball-in-a-cup task, the robotic arm is mounted on a stand as shown in Figure 1 (right). The Ball-in-a-cup cell is equipped with two cameras positioned orthogonally – and both facing the robot – as well as with a Vicon Vero setup. The Vicon system is solely used for computing a sparse external "catch" reward.

In all experiments we use Scheduled Auxiliary Control (SAC-X) with a scheduler choosing a sequence of 3 intentions per episode. Since the environment is initialized such that the initial responses are uniformly distributed, the expected return of the 'increase' and 'decrease' rewards following the optimal policy $\pi_k^*$ for a respective sensor reward is

$$V^{\pi_k^*}(s) \leq \frac{|z_{\max}^k - z_{\min}^k|}{2} \quad \forall\, s \in \mathcal{S}.$$

Accordingly, we scale the increase and decrease rewards by a constant factor $\frac{2\sigma}{|z_{\max} - z_{\min}|}$ in all experiments and choose $\sigma = 200$ (which corresponds to the number of steps an SSI is executed in an episode) to achieve comparable reward scales across the reward schemes.

In all simulated setups we use uniform scheduling (SAC-U), where intentions are chosen by uniform random sampling [33]. Also, we use a multi-actor setup with 64 concurrent actors that collect data in parallel. This approach optimizes for wall-clock time and sacrifices data-efficiency, which we accepted, since these experiments were primarily intended to answer the general question whether SSIs allow us to learn sparse tasks. The question of data-efficiency is relevant, however, in the real robot experiments and we thus used a jointly learned scheduler to select between policies to execute (SAC-Q, [33]) for improved data-efficiency. All real world experiments were conducted using a single robot (corresponding to a single-actor setup in simulation).
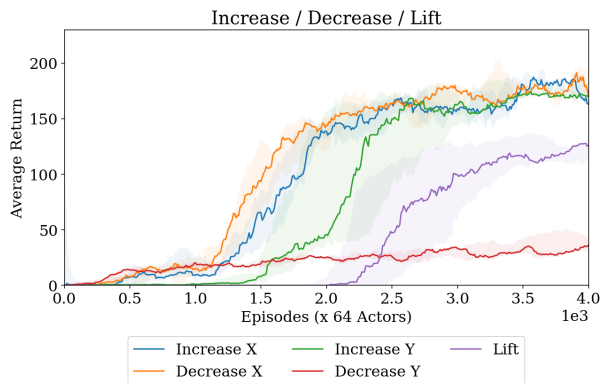
Fig. 3. 'Lift' learned from pixels in the simulated manipulation setup with the 'increase' and 'decrease' rewards used as auxiliary intentions.



Fig. 4. Comparison of the different reward schemes for learning 'Lift' from pixels in the simulated manipulation setup.

In all cases, the agent is provided with an observation that comprises proprioceptive information - joint positions, velocities and torques - as well as a wrist sensor's force and torque readings in the manipulation setups. Additionally, the agent receives the camera images as exteroceptive inputs. Thus, all tasks have to be learned from proprioceptive sensors and raw pixel inputs. The action space in the manipulation setups is five dimensional and continuous and consists of the three Cartesian translational velocities, the angular velocity of the wrist around the vertical axis and the speed of the grippers fingers. The action space in the Ball-in-a-Cup setup is four dimensional and continuous and consists of the raw target joint velocities [38]. In all cases, the robot is controlled at 20 Hz.

### B. Learning to grasp and lift in simulation

Grasping and lifting an object with a robotic arm is a challenging task, in particular when learned from scratch and purely from pixels: The agent must learn to recognize the object, approach it, find the right position for grasping and eventually close the fingers and lift the object. Learning this from scratch, when only a sparse final reward is given, is very unlikely. We assume the external target reward for lift is given by a binary signal: the external reward is 1, if the touch sensor is triggered and the gripper is at least 15cm above the basket.

In a first experiment, we are using the 'delta response rewards' SSIs, that give reward for pushing the distribution of pixels of a selected color channel in a certain direction. In this experiment, we select the color channel to roughly match the color of the block. However, we will discuss below in section VI-C, how we can get rid of this assumption, or make it more general, respectively. This selection of SSIs results in overall four auxiliary intentions, two for increasing the mean in x- or y-direction, and two for decreasing the mean in x- or y-direction of the pixel image. The learning curves for the four auxiliaries and the final lift reward (violet line) are shown in Figure 3. After about 500 episodes (times 64 actors), the agent sees some reward for moving the mean in various directions. This results in first, small interactions of the arm with the object. After about 1000 episodes, interactions get
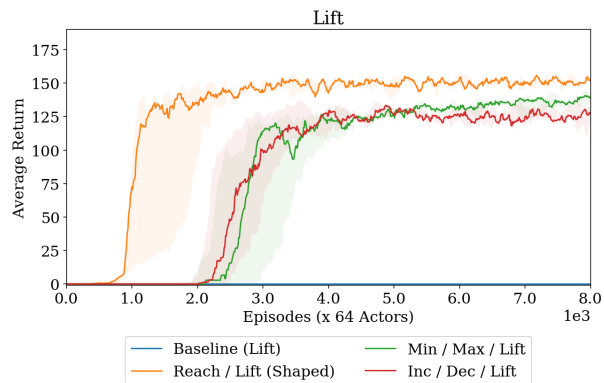
stronger, until after about 2000 episodes, the agent learns to lift the object deliberately. Final performance is reached after about 4000 episodes. We note that the reward curve for the 'decrease y' intention is lower than the other curves because of the fact, that the block is typically at the bottom of the basket and therefore the mean of the color channel's distribution is typically already pretty low, so 'decrease y' is not able to earn as much reward as the other auxiliaries. This is an example of an intention that is potentially not very useful, but however does not prevent the agent of finally learning the goal task.

If no auxiliary rewards are given, the agent does not learn to lift the object at all. This is shown by the flat blue learning curve in Figure 4. This Figure also shows the learning curve for the alternative 'target response reward' SSI formulation, which tries to minimize/ maximize the value of the mean (green line) in comparison to the 'delta response reward' SSI formulation (red line). For the lift experiment, not much difference in the learning curves can be seen; both SSI formulations work successfully in this setting. For further reference, we also conducted a typical learning experiment using a dedicated shaping reward for reaching and grasping the block (orange line). As expected, the shaped approach learns faster, but at the costs of a considerable effort in the specification of the shaping reward and the required instrumentation for computing the position the object (which would require camera calibration, object detection, and pose estimation on a real system).

### C. Ablation studies

We conducted several ablation studies and summarize our findings below:

*1) Learning success does not depend on a particular camera pose for reward giving:* We investigated this by varying the perspective of the camera, that is used for reward computation, in various ways (see Figure 5). As shown by the learning curves in Figure 6, the concrete pose of the camera has an influence on learning speed, but in all cases, the SSIs were powerful enough to learn the final lifting task eventually.
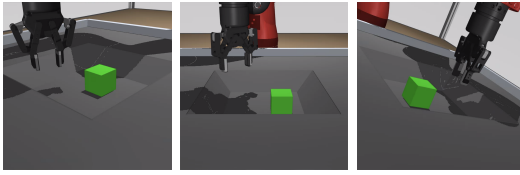
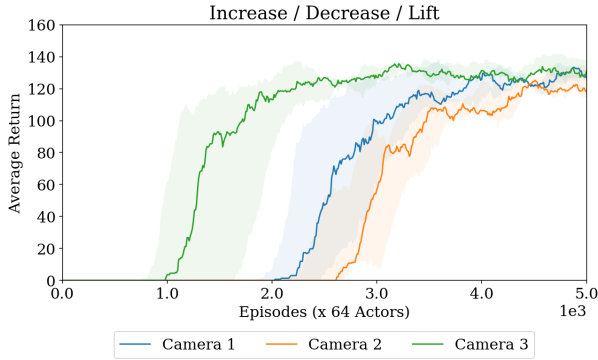Fig. 5. The three camera angles used for computing the SSIs for the experiments shown in Figure 6.



Fig. 6. Comparison of different camera angles for computing SSIs when learning the 'Lift' task. In all experiments the 'increase' and 'decrease' rewards are used as auxiliary intentions.

*2) The SSI color channel does not necessarily need to specify a single object:* We investigated this by two experiments: an experiment with two blocks of the same color (see figure 7, left) and an experiment, where part of the background had the same color as the block (see Figure 7, right). In both cases lift could be learned successfully using the standard pixel based SSIs. If we increase the proportion of the background pixels with the same color as the object, at some point the agent fails to learn to interact with the brick, but instead 'exploits' the reward scheme by trying to move the mean by hiding pixels with the arm and gripper. It is possible to filter the background from the scene; however, this is beyond the scope of this paper and therefore left for future investigations.

*3) One can use a much more general selection of the color channel:* To demonstrate this, we used the 'aggregate' SSI formulation suggested above: We compute rewards (here: delta rewards) for a potentially large set of different color channels and add those up in one so-called 'aggregate' SSI. This will clearly solve the single block lift task as described above, as long as at least one of the color channels matches the
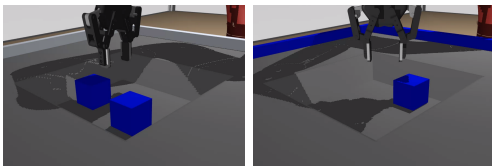


Fig. 7. The setups used for demonstrating robustness. On the left, both objects in the scene have the same color. On the right a non-moveable part of the basket has the same color as the object.
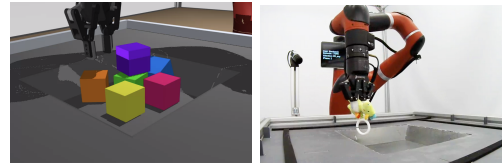


Fig. 8. Left: A subset of the colored blocks used in the simulated 'Lift Any' experiments. In each episode one of the available blocks is selected. Right: Evaluation of the 'Lift Any' experiment showing that the robot is able to grasp and lift a non-rigid, multi-colored baby toy.
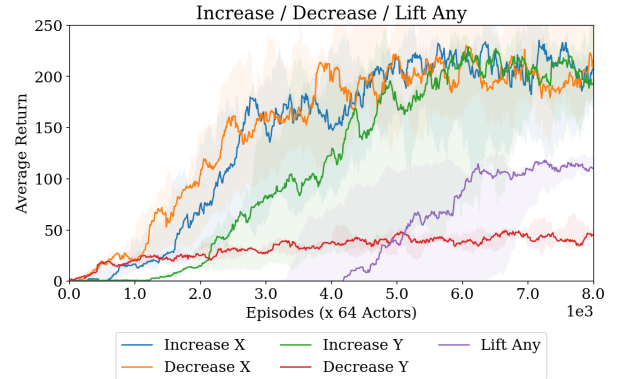


Fig. 9. 'Lift Any' learned from pixels in the simulated manipulation setup with the 'increase' and 'decrease' rewards used as auxiliary intentions. Here, the 'increase' and 'decrease' rewards are aggregated over color channels and the block's color is changed every episode.

color of the block. However, this works also in a setting, where blocks with different colors are used throughout the experiment. Figure 9 shows the learning curve for the 'Lift Any' task (red line), where we randomly changed the color of the block in every episode (see Figure 8 showing a subset of the colored blocks used in the experiment). A similar experiment was also conducted on the real robot. The results for the real world experiment is shown in Figure 12.

*4) The SSI method is not restricted to pixels, but works with a general set of (robot) sensors:* In particular, we conducted experiments in a setup with two blocks in the scene, where we applied SSIs to basic sensors, like the touch sensor and the joint angles, but also used the SSIs described before on camera images, resulting in a total of 22 auxiliary intentions (minimize / maximize touch, minimize / maximize joint angles of arm and gripper joints, minimize / maximize the mean of the color distribution along the x and y axes). To deal efficiently with the extended set of intentions, the agent used a Q-based scheduler (SAC-Q), which successfully learned the 'Lift' task.

*5) It is mandatory to have a penalty term for moving the sensor response in the opposite direction:* In learning experiments, where the agent only received a positive reward, if the response is moved in the intended direction, the agent quickly learns to cleverly exploit the reward, e.g. by moving the gripper back and forth in front of camera, hiding and revealing the block and thereby collecting reward for moving the response (see III-A and III-B).
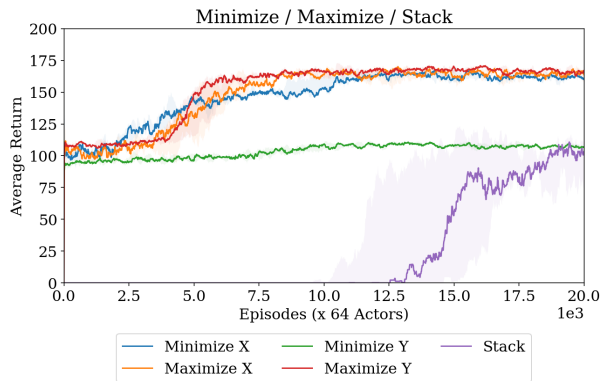
Fig. 10. 'Stack' learned from pixels in the simulated manipulation setup with the 'minimize' and 'maximize' rewards used as auxiliary intentions.



Fig. 11. 'Lift' learned from pixels on the real robot with the 'increase' and 'decrease' rewards used as auxiliary intentions. In this setup, the 'increase' and 'decrease' rewards are aggregated over two perpendicular camera angles.

### D. Learning to stack

Learning to stack a block on another block poses additional challenges compared to the 'grasp-and-lift' scenario described above: The scene is more complex since there are two objects now, reward is given only if the object is placed above the target object, and the target object can move. The external task reward for 'stack' is 1, if and only if the block 1 is properly placed on block 2 and the gripper is open.

We use the minimize/ maximize SSI approach in combination with a SAC-U agent: the SSI auxiliary rewards are computed from raw pixels and the input observations for the agent comprises raw pixels and proprioceptive information only. Figure 10 shows the learning curves for the auxiliaries and the final 'Stack' reward. This is a much more difficult task to learn from pure pixels, but after about 20,000 episodes (times 64 actors) the agent has figured out how to reliably solve the task, purely from pixels. Not surprisingly, without SSI auxiliaries, the agent is not able to learn the task. The minimize / maximize auxiliary rewards are position based rewards - therefore the learning curves are offset by roughly the average reward of 100 which corresponds to the usual average position of the pixel mean.

Although the amount of data used to learn this task (20,000 episodes times 64 actors) is huge, it is still surprising, that learning such a complex task is possible from raw sensor information and an external task reward only. We assume this is possible, since the used simple sensor intentions encourage a rich playing within the scene, and the additional sequencing of these intentions as performed by SAC-X increases the probability of seeing also more complex configurations - like stacking in this case. From seeing external rewards for these occasional stacks, learning finally can take off.

### E. Learning to grasp and lift on a real robot

To investigate the behaviour of SSI based agents in a real world robotic setup, we apply the agent to learn to grasp and lift objects (figure 1, left). As in the simulation experiments, the agent's input is based on raw proprioceptive sensory information 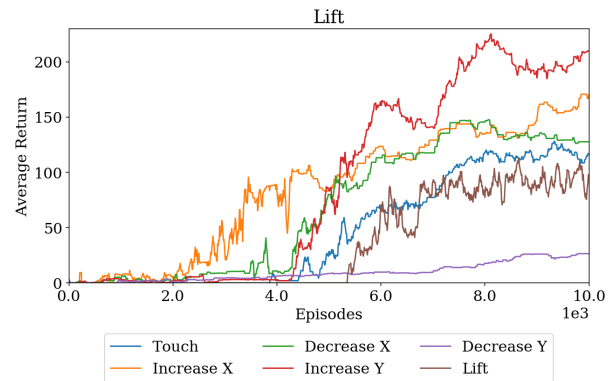and the raw images of two cameras placed around the basket. Real world experiments add the additional challenge of noisiness of sensors and actuators, with which the agent has to cope. Also, the approach naturally has to work in a single actor setup, since there is only one robot available to collect the data.

For the real robot experiment, we use 6 SSIs based on the touch sensor as well as on the camera images: increase/ decrease of the color distribution's mean in x- and y-direction of raw camera images (4 intentions) plus minimize / maximize touch sensor value ('on' / 'off', 2 intentions). The task reward for 'Lift' is given sparsely, if and only if the touch sensor is activated while the gripper is 15 cm above the table. To make the SSI reward signal from pixels more robust, we first computed the SSI values from each camera image and then aggregated the rewards to get one single accumulated reward. This means that an auxiliary task receives the maximum reward only if it achieves to move the color distribution's mean in both cameras in the desired direction.

The learning agent is again based on Scheduled Auxiliary Control and we apply a Q-table based scheduler (SAC-Q) [33] to make learning as data-efficient as possible. We also add a special exploration scheme based on 'multi-step' actions [36, 37]: when selecting an action, we additionally determine how often to repeat its application, by drawing a number of action repeats uniformly in the range $[1, 20]$. Also, for safety reasons, external forces are measured at the wrist sensor and the episode is terminated if a threshold of 20N on any of the three principle axes is exceeded for more than three time steps in a row. If this happens, the episode is terminated with zero reward, enforcing the agent to avoid these failure situations in the future.

We find that the agent successfully learns to lift the block as illustrated in Figure 11. After an initial phase of playful, pushing-style interaction, the agent discovers a policy for moving the block to the sides of the basket and after exploring the first successful touch rewards, it quickly learns to grasp and lift the block. After about 9000 episodes the agent shows a reliable lifting behaviour. This corresponds to roughly 6 days of training on the real robot.
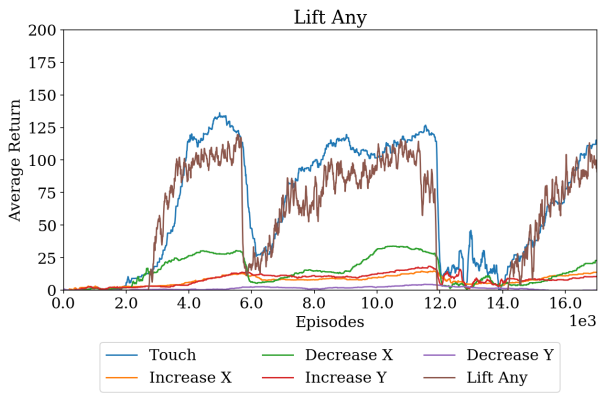
Fig. 12. 'Lift Any' learned from pixels on a real robot with the 'increase' and 'decrease' rewards used as auxiliary intentions. The block is replaced at various points throughout the experiment.



Fig. 13. 'Catch' learned from pixels on a real robot.

In a further experiment, we replaced the above SSIs for a single color channel with SSIs, that aggregate rewards over multiple color channels, allowing to learn with objects of any color. We tested this by starting to learn with a single object, until the robot started to lift, and then replacing the object by another object of different color and/ or of different shape. The learning curve is shown in figure 12. The drops in the learning curve indicate, that if the object is replaced, the agent does not know how to lift yet. After some time of adaptation, it starts to manage to lift again. Continuing this, we saw the robot being able to learn to lift a wide variety of different objects, all from raw sensor information.

### F. Ball in a Cup

An important aspect of simple sensor intentions is, that SSIs constitute a basic concept and show some generality of being helpful for different external target tasks. To illustrate how the same set of SSIs can be employed to master a completely different control problem, we show results on the dynamic Ball-in-a-Cup task [38]: the task is to swing up a ball attached to a robot arm and to catch it with a cup mounted as the end effector of the arm. The agent only receives a binary positive reward, if the ball is caught in the cup (see figure 1, right).

Dynamic tasks in general exhibit additional difficulties compared to static tasks, e.g. the importance of timing and the difficulty of reaching (and staying in) possibly unstable regimes of the robot's configuration-space. As a result, learning to catch the ball purely from pixels is out-of-reach for learning setups, that only employ the sparse catch reward.

To show the versatility of simple sensor intentions, we choose the standard increase / decrease SSIs, resulting in 4 auxiliary tasks plus the binary task reward for catch. The learning agent used is a SAC-Q agent. As shown in Figure 13, the agent is able to learn this dynamic task purely from the standard observation, pixels and proprioceptive inputs. In order to cope with the dynamic nature of the task, 2 consecutive pixel frames were stacked for the controller's input.

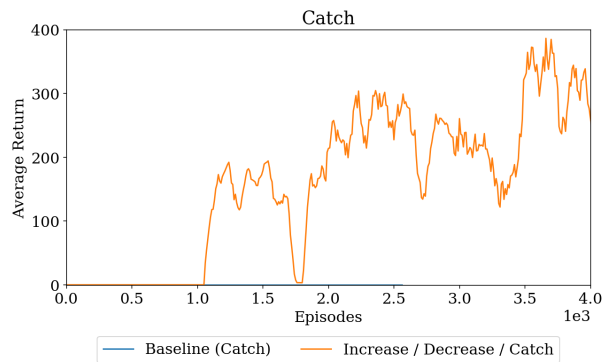The simple sensor intentions encourage to move the colored

pixels in the image, which results in learning to deliberately move the ball. The sequential combination of different auxiliary tasks as enforced by the SAC-X agent, lead to a broad exploration of different movements of ball and robot. Eventually, first catches are observed, and once this happens, the agent quickly learns the 'Catch' task. For the roughly 4000 episodes needed, the learning took about 3 days in real time.

## VII. RELATED WORK

Transfer via additional tasks has a long-standing history in reinforcement learning to address exploration challenges and accelerate learning [45, 31]. We are able to conceptually distinguish into two main categories with auxiliary tasks [43, 23], which are used to accelerate training on a final task, and multitask learning [9] with focus on the performance of all involved tasks.

Early work on multitask learning [12] introduced the prediction of expected sums of future values for multiple tasks of interest. Many successive approaches have extended modelling of multiple task rewards and investigated further types of transfer across tasks [42, 34, 4]. Similarly, work on the options framework investigates directions for decomposition of task solutions [13, 3, 11, 48].

Auxiliary tasks have been investigated as manually chosen to help in specific domains [33, 14, 23, 26, 8] and as based on agent behaviour [2]. As manual task design presents significant burden for human operators, these methods demonstrate that even simpler sets of tasks can demonstrate considerable benefits. In comparison to methods using auxiliary tasks mostly for representation shaping by sharing a subset of network parameters across tasks [23, 26], SSI shares data between tasks which directly uses additional tasks for exploration.

For both directions, success considerably relies on the source of tasks and correlated reward functions, which this work focuses on. Multiple agent behaviour based sources for tasks and additional objectives have been considered in prior work including diversity of behaviours [39, 19, 15], intrinsic motivation [10, 41, 7, 30, 25], and empowerment of the agent [24, 28, 22]. Additionally, representation learning and the automatic identification of independently controllable features have provided another perspective on identifying tasks for

transfer and improving exploration [19, 5, 7]. Recent work on diversity has in particular demonstrated the importance of the space used for skill discovery [39]. SSI provides a valuable perspective on determining valuable task spaces with limited human effort.

Given large sets of tasks, the higher-level problem of choosing which task to learn creates a novel challenge similar to exploration within a single task. Often random sampling over all possible tasks provides a strong baseline [17]. However, to accelerate learning different perspectives build on curriculum [6, 21, 16], iterative task generation [35, 47]. In this work, we rely on task scheduling similar to Riedmiller et al. [33] in order to optimize the use of training time.

## VIII. Conclusion and future work

Learning to change sensor responses deliberately is a promising exploration principle in settings, where it is difficult or impossible to experience an external task reward purely by chance. We introduce the concept of simple sensor intentions (SSIs) that implements the above principle in a generic way within the SAC-X framework. While the general concept of SSIs applies to any robotic sensor, the application to more complex sensors, like camera images is not straight forward. We provide one concrete way to implement the SSI idea for camera images, which first need to be mapped to scalar values to fit into the proposed reward scheme. We argue that our approach requires less prior knowledge than the broadly used shaping reward formulation, that typically rely on task insight for their definition and state estimation for their computation.

In several case studies we demonstrated the successful application to various robotic domains, both in simulation and on real robots. The SSIs we experimented with were mostly based on pixels, but also touch and joint angle based SSIs were used. The definition of the SSIs was straight-forward with no or minor adaptation between domains.

Future work will concentrate on the extension of this concept in various directions, e.g. improving the scheduling of intentions to deal with a large number of auxiliary tasks will enable the automatic generation of rewards and reward combinations.

## Acknowledgments

## References

[1] Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Rémi Munos, Nicolas Heess, and Martin A. Riedmiller. Maximum a posteriori policy optimisation. *CoRR*, abs/1806.06920, 2018.

[2] Marcin Andrychowicz, Dwight Crow, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pages 5055–5065, 2017.

[3] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

[4] Andre Barreto, Will Dabney, Remi Munos, Jonathan J Hunt, Tom Schaul, David Silver, and Hado P van Hasselt. Successor Features for transfer in reinforcement learning. In *Advances in Neural Information Processing Systems 30 (NIPS)*. 2017.

[5] Emmanuel Bengio, Valentin Thomas, Joelle Pineau, Doina Precup, and Yoshua Bengio. Independently controllable features. *arXiv preprint arXiv:1703.07718*, 2017.

[6] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, 2009.

[7] Sebastian Blaes, Marin Vlastelica Pogančić, Jiajie Zhu, and Georg Martius. Control what you can: Intrinsically motivated task-planning agent. In *Advances in Neural Information Processing Systems*, pages 12520–12531, 2019.

[8] Serkan Cabi, Sergio Gomez Colmenarejo, Matthew W. Hoffman, Misha Denil, Ziyu Wang, and Nando de Freitas. The Intentional Unintentional Agent: Learning to solve many continuous control tasks simultaneously. In *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, Proceedings*, pages 207–216, 2017.

[9] Rich Caruana. Multitask learning. *Machine learning*, 28 (1):41–75, 1997.

[10] Nuttapong Chentanez, Andrew G. Barto, and Satinder P. Singh. Intrinsically Motivated reinforcement learning. In *Advances in Neural Information Processing Systems 17 (NIPS)*. 2005.

[11] Christian Daniel, Gerhard Neumann, Oliver Kroemer, and Jan Peters. Hierarchical relative entropy policy search. *The Journal of Machine Learning Research*, 17 (1):3190–3239, 2016.

[12] Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 1993.

[13] Thomas G. Dietterich. The maxq method for hierarchical reinforcement learning. In *In Proceedings of the Fifteenth International Conference on Machine Learning (ICML)*, 1998.

[14] Alexey Dosovitskiy and Vladlen Koltun. Learning to act by predicting the future. In *International Conference on Learning Representations (ICLR)*, 2017.

[15] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.

[16] Sébastien Forestier, Yoan Mollard, and Pierre-Yves

Oudeyer. Intrinsically motivated goal exploration processes with automatic curriculum learning. *arXiv preprint arXiv:1708.02190*, 2017.

[17] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1311–1320. JMLR. org, 2017.

[18] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. In *International Conference on Learning Representations*, 2017.

[19] Christopher Grimm, Irina Higgins, Andre Barreto, Denis Teplyashin, Markus Wulfmeier, Tim Hertweck, Raia Hadsell, and Satinder Singh. Disentangled cumulants help successor representations transfer to new tasks. *arXiv preprint arXiv:1911.10866*, 2019.

[20] Nicolas Heess, Gregory Wayne, David Silver, Tim Lillicrap, Tom Erez, and Yuval Tassa. Learning continuous control policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*, pages 2944–2952, 2015.

[21] Nicolas Heess, Dhruva Tirumala, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, Ali Eslami, Martin Riedmiller, et al. Emergence of locomotion behaviours in rich environments. *arXiv preprint arXiv:1707.02286*, 2017.

[22] Rein Houthooft, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems*, pages 1109–1117, 2016.

[23] Max Jaderberg, V Mnih, W M Czarnecki, T Schaul, J Z Leibo, D Silver, and K Kavukcuoglu. Unreal: Reinforcement learning with unsupervised auxiliary tasks. In *ICLR 2017*, 2017.

[24] Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. All else being equal be empowered. In *European Conference on Artificial Life*, pages 744–753. Springer, 2005.

[25] Tejas D. Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J. Gershman. Deep Successor reinforcement learning. *arXiv preprint arXiv:1606.02396*, 2016.

[26] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to Navigate in complex environments. *arXiv preprint arXiv:1611.03673*, 2016.

[27] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015. ISSN 00280836. URL http://dx.doi.org/10.1038/nature14236.

[28] Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. In *Advances in neural information processing systems*, pages 2125–2133, 2015.

[29] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287, 1999.

[30] Hung Quoc Ngo, Matthew D. Luciw, Alexander Förster, and Jürgen Schmidhuber. Learning skills from play: Artificial curiosity on a katana robot arm. In *The 2012 International Joint Conference on Neural Networks (IJCNN)*, 2012.

[31] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.

[32] Jette Randløv and Preben Alstrøm. Learning to drive a bicycle using reinforcement learning and shaping. In *ICML*, volume 98, pages 463–471. Citeseer, 1998.

[33] Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degrave, Tom Wiele, Vlad Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing solving sparse reward tasks from scratch. In *International Conference on Machine Learning*, pages 4341–4350, 2018.

[34] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal Value function approximators. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.

[35] Jürgen Schmidhuber. PowerPlay: Training an increasingly general problem solver by continually searching for the simplest still unsolvable problem. In *Front. Psychol.*, 2013.

[36] R. Schoknecht and M. Riedmiller. Speeding up reinforcement learning with multi-step actions. In *Proceedings of International Conference on Artificial Neuraln Networks, ICANN'02*, Madrid, Spain, 2002.

[37] R. Schoknecht and M. Riedmiller. Reinforcement learning on explicitly specified time-scales. *Neural Computing & Applications*, 12(2):61–80, November 2003.

[38] Devin Schwab, Jost Tobias Springenberg, Murilo Fernandes Martins, Michael Neunert, Thomas Lampe, Abbas Abdolmaleki, Tim Hertweck, Roland Hafner, Francesco Nori, and Martin A. Riedmiller. Simultaneously learning vision and feature-based control policies for real-world ball-in-a-cup. In Antonio Bicchi, Hadas Kress-Gazit, and Seth Hutchinson, editors, *Robotics: Science and Systems XV, University of Freiburg, Freiburg im Breisgau, Germany, June 22-26, 2019*, 2019. doi: 10.15607/RSS. 2019.XV.027. URL https://doi.org/10.15607/RSS.2019. XV.027.

[39] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019.

[40] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, jan 2016. ISSN 0028-0836. doi: 10.1038/nature16961.

[41] S. Singh, R. L. Lewis, and A. G. Barto. Where do rewards come from? In N.A. Taatgen and H. van Rijn, editors, *Proceedings of the 31st Annual Conference of the Cognitive Science Society*. Austin, TX, 2009.

[42] Richard S. Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M. Pilarski, Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems - Volume 2*, AAMAS '11, pages 761–768, Richland, SC, 2011. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 0-9826571-6-1, 978-0-9826571-6-4.

[43] Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(Jul):1633–1685, 2009.

[44] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 5026–5033. IEEE, 2012.

[45] Lisa Torrey and Jude Shavlik. Transfer learning. In *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*, pages 242–264. IGI Global, 2010.

[46] Vicon Motion Systems Ltd. U.K. Tracker motion capture software for vr object tracking, 2019. URL https://www.vicon.com/products/software/tracker. [Online; accessed 14-April-2020].

[47] Rui Wang, Joel Lehman, Jeff Clune, and Kenneth O Stanley. Paired open-ended trailblazer (poet): Endlessly generating increasingly complex and diverse learning environments and their solutions. *arXiv preprint arXiv:1901.01753*, 2019.

[48] Markus Wulfmeier, Abbas Abdolmaleki, Roland Hafner, Jost Tobias Springenberg, Michael Neunert, Tim Hertweck, Thomas Lampe, Noah Siegel, Nicolas Heess, and Martin Riedmiller. Regularized hierarchical policies for compositional transfer in robotics. *arXiv preprint arXiv:1906.11228*, 2019.

In the following sections we give a detailed description of the experimental setup for the manipulation and the Ball-in-a-Cup domains.

### A. Manipulation setup

For the experiments in the manipulation setup we use a MuJoCo [44] simulation that is well aligned with the real world robot's setup. In both cases, a Rethink Sawyer robotic arm, with a Robotiq 2F-85 parallel gripper as end-effector, faces a 20cm x 20xm basket, that, in the case of the 'Lift' and 'Lift Any' tasks, contains a single colored block or, in the case of the 'Stack' task, two differently colored blocks. The basket is equipped with three cameras, two attached at the front corners of the basket and one attached in the back-left corner. All cameras face the basket's center. The 3D-printed, colored blocks are cubic, with a side length of 5cm.

In simulation, the blocks are initialized at random locations within the 20cm x 20cm workspace at the beginning of each episode. On the real robot an initialization sequence, that randomizes the block positions, is triggered every 20 episodes.

As shown in Table I, the agent receives proprioceptive information - the joint positions, joint velocities and joint torques, as well as the wrist sensor's force torque readings. Additionally, the agent is provided with a binary grasp sensor. Further, the agent receives the camera images of all three cameras that are attached to the basket as the only exteroceptive inputs. As shown in Table II, the action space is five dimensional, continuous and consists of the three cartesian translational velocities, the angular velocity of the wrist around the vertical axis and the current speed of the gripper's fingers. The workspace is limited to the 20cm$^3$ cube above the table surface.

To determine a successful lift, we compute the end-effector's tool center point (TCP) using forward kinematics. The 'Lift' and 'Lift Any' rewards are given in terms of the TCP's z-coordinate $\text{TCP}_z$ by

$$r_{\text{lift}} = \begin{cases} 1 & \text{if } \text{TCP}_z > 0.15\text{m and the grasp sensor is active,} \\ 0 & \text{otherwise.} \end{cases}$$

Additionally, on the real robot, we measure external forces at the wrist sensor for safety reasons and the episode is terminated if a threshold of 20N on any of the three principle axes is exceeded for more than three consecutive time steps. If this happens, the episode is terminated with zero reward, enforcing the agent to avoid these failure situations in the future.

### B. Ball-in-a-Cup

The Ball-in-a-Cup setup features a Rethink Sawyer robotic arm mounted on a stand, with a custom, 3D-printed cup attachment, which is fixed to the arm's wrist. The cup has a diameter of 20cm and a height of 16cm. A woven net is attached to the cup, to ensure, that the ball is visible even

when it is inside the cup. The ball has a diameter of 5cm and is attached to a ball bearing at the robot's wrist with 40cm Kevlar string. The bearing helps to prevent winding of the string – when the ball goes around the wrist, the string rotates freely via the bearing. The Ball-in-a-cup cell is equipped with two cameras positioned orthogonally – and both facing the robot – as well as with a Vicon Vero [46] setup, that tracks the cup and the ball and which is used for computing the sparse, external 'catch' reward, as well as for triggering resets.

TABLE I
OBSERVATION SPACE OF THE MANIPULATION SETUP.

| Entry | Dimensions | Unit |
|---|---|---|
| Joint Position (Arm) | 7 | rad |
| Joint Velocity (Arm) | 7 | rad/s |
| Joint Torque (Arm) | 7 | Nm |
| Joint Position (Hand) | 1 | rad |
| Joint Velocity (Hand) | 1 | tics/s |
| Force-Torque (Wrist) | 6 | N, Nm |
| Binary Grasp Sensor | 1 | au |
| TCP Pose | 7 | m, au |
| Camera Image (Front Right) | 64 x 64 x 3 | |
| Camera Image (Front Left) | 64 x 64 x 3 | |
| Camera Image (Back Left) | 64 x 64 x 3 | |
| Last Control Command (Joint Velocity) | 8 | rad/s |

TABLE II
ACTION SPACE OF THE MANIPULATION SETUP.

| Entry | Dimensions | Unit | Range |
|---|---|---|---|
| Translational Velocity (in x, y, z) | 3 | m/s | [-0.07, 0.07] |
| Wrist Rotation Velocity | 1 | rad/s | [-1, 1] |
| Finger speed | 1 | tics/s | [-255, 255] |

TABLE III
OBSERVATION SPACE OF THE BALL-IN-A-CUP SETUP.

| Entry | Dimensions | Unit |
|---|---|---|
| Joint Position (Arm) | 7 | rad |
| Joint Velocity (Arm) | 7 | rad/s |
| Camera Image (Frontal) | 84 x 84 x 3 | |
| Camera Image (Side) | 84 x 84 x 3 | |
| Last Control Command (Joint Velocity) | 4 | rad/s |
| Action Filter State | 4 | rad/s |

TABLE IV
JOINT LIMITS OF THE BALL-IN-A-CUP SETUP.

| | Unit | J0 | J1 | J5 | J6 |
|---|---|---|---|---|---|
| Position | rad | [-0.4, 0.4] | [0.3, 0.8] | [0.5, 1.34] | [2.6, 4.0] |
| Velocity | rad/s | [-2.0, 2.0] | [-2.0, 2.0] | [-2.0, 2.0] | [-2.0, 2.0] |

TABLE V
HYPERPARAMETERS FOR MPO AND SAC-X.

| | |
|---|---|
| 2D Conv Channels | 128, 64, 64 |
| 2D Conv Kernel Sizes | 4x4, 3x3, 3x3 |
| 2D Conv Strides | 2, 2, 2 |
| Actor Torso Layer Sizes | 256 |
| Actor Head Layer Sizes | 100, 16 or 8 |
| Critic Torso Layer Sizes | 400 |
| Critic Head Layer Sizes | 300 , 1 |
| Activation Function | elu |
| Discount Factor | 0.99 |
| Adam Learning Rate | 2e-4 |
| Replay Buffer Size | 1e6 |
| Target Network Update Period | 500 |
| Batch Size | 64 |
| KL Constraint Epsilon | 1 |
| KL Epsilon Samples | 0.1 |
| KL Epsilon Mean | 1e-3 |
| KL Epsilon Covariance | 1e-5 |

gradients [20] are used for optimization, we use Maximum a Posteriori Policy Optimisation (MPO) [1] for learning the policy and value function. Following [33], we implement both the policy and value function with multi-headed (multi-task) neural networks. All pixel inputs are passed through a convolutional neural network with three convolutional layers. The resulting embeddings are subsequently concatenated with the proprioceptive inputs. The torso is shared between the network heads and comprises a single fully-connected layer. The task heads in turn consist of two additional fully-connected layers. The final layer outputs the parameters of a multivariate Gaussian distribution, with the same dimensionality as the action space. The network details and hyperparameter values are shown in Table V.

The action space 4 dimensional, continuous and consists of the joint velocities to be commanded for 4 of the 7 joints (J0, J1, J5 and J6), which are sufficient to solve the task. The unused degrees of freedom are commanded to remain at a fixed position throughout the episode. The imposed joint limits are shown in Table IV and are chosen as to avoid collisions with the workspace and to keep the ball and cup visible in both cameras at all times. To prevent high frequency oscillations in the commanded velocities, the outputs of the policy network are passed through a low-pass filter with a cutoff frequency of 0.5 Hz.

As shown in Table III, the agent receives proprioceptive information - the joint positions and joint velocities, as well as the camera images of the two cameras attached to the cell. Further, the agent is provided with the action executed at the previous time step as well as with the low-pass filter's internal state.

Due to the string wrapping and tangling around the robot and cup, resets can be challenging. In order to minimize human intervention during training, the robot uses a simple hand-coded reset policy that allows for mostly unattended training. We refer to [38] for details on the reset procedure.

We use the 3D pose of the ball and the cup, obtained from the Vicon Vero tracking system, to determine the sparse, external reward, which is only given in case of a successful catch:

$$r_{\text{catch}} = \begin{cases} 1 & \text{if ball in cup,} \\ 0 & \text{otherwise.} \end{cases}$$

## APPENDIX B
## NETWORK ARCHITECTURE AND HYPERPARAMETERS

In all experiments we use Scheduled Auxiliary Control (SAC-X) [33], either with a uniform random scheduler (SAC-U) or with a learned, Q-table based scheduler (SAC-Q). However, in contrast to the original work, where stochastic value