

---

# SOLVING RUBIK’S CUBE WITH A ROBOT HAND

---

A PREPRINT

**OpenAI**

Ilge Akkaya\*, Marcin Andrychowicz\*, Maciek Chociej\*, Mateusz Litwin\*, Bob McGrew\*, Arthur Petron\*,  
Alex Paino\*, Matthias Plappert\*, Glenn Powell\*, Raphael Ribas\*, Jonas Schneider\*, Nikolas Tezak\*,  
Jerry Tworek\*, Peter Welinder\*, Lilian Weng\*, Qiming Yuan\*, Wojciech Zaremba\*, Lei Zhang\*

October 14, 2019

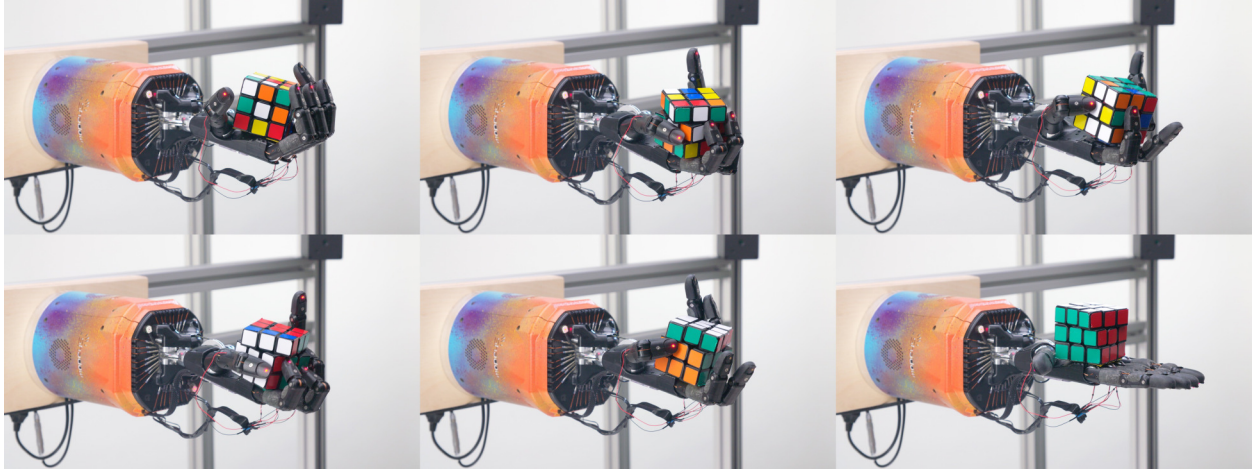


Figure 1: A five-fingered humanoid hand trained with reinforcement learning and automatic domain randomization solving a Rubik’s cube.

## ABSTRACT

We demonstrate that models trained only in simulation can be used to solve a manipulation problem of unprecedented complexity on a real robot. This is made possible by two key components: a novel algorithm, which we call automatic domain randomization (ADR) and a robot platform built for machine learning. ADR automatically generates a distribution over randomized environments of ever-increasing difficulty. Control policies and vision state estimators trained with ADR exhibit vastly improved sim2real transfer. For control policies, memory-augmented models trained on an ADR-generated distribution of environments show clear signs of emergent meta-learning at test time. The combination of ADR with our custom robot platform allows us to solve a Rubik’s cube with a humanoid robot hand, which involves both control and state estimation problems. Videos summarizing our results are available: <https://openai.com/blog/solving-rubiks-cube/>

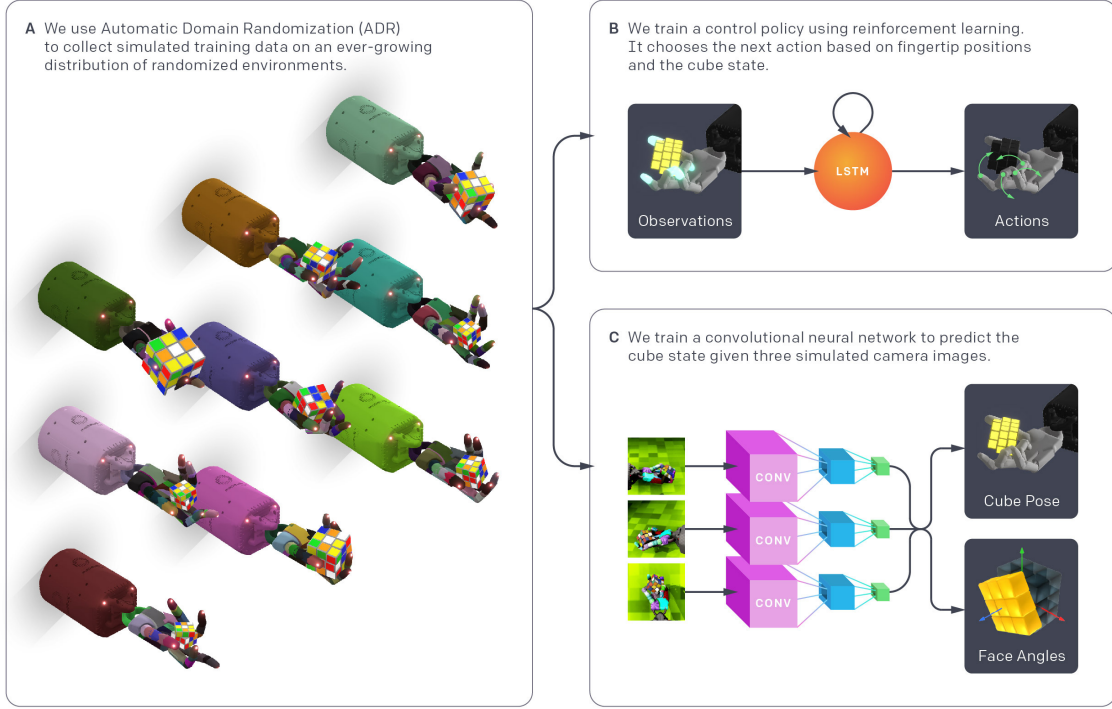
## 1 Introduction

Building robots that are as versatile as humans remains a grand challenge of robotics. While humanoid robotics systems exist [1, 100, 5, 96, 10], using them in the real world for complex tasks remains a daunting challenge. Machine learning

---

\* Authors are listed alphabetically. We include a detailed contribution statement at the end of this manuscript. Please cite as OpenAI et al., and use the following bibtex for citation: <https://openai.com/bibtex/openai2019rubiks.bib>

## Train in Simulation



## Transfer to the Real World

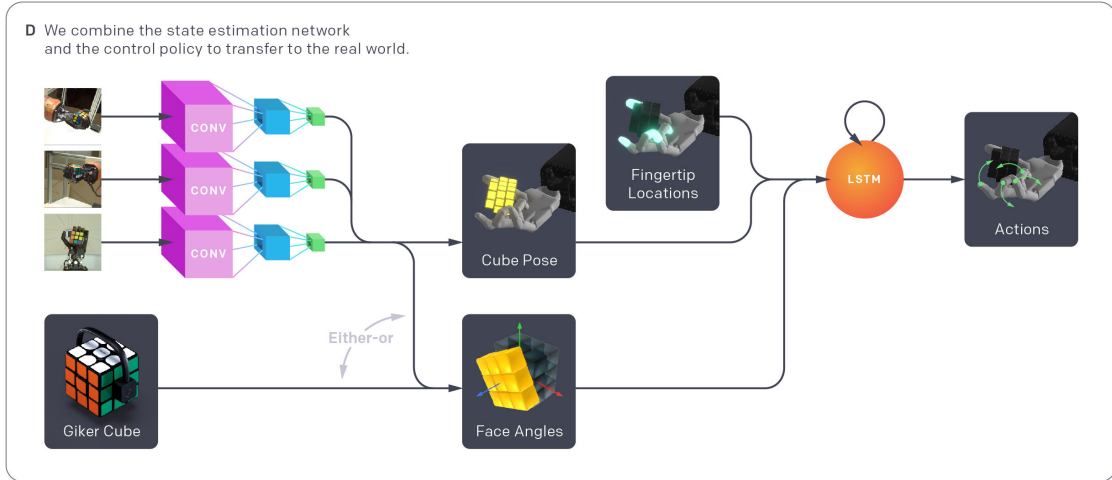


Figure 2: System Overview. (a) We use automatic domain randomization (ADR) to generate a growing distribution of simulations with randomized parameters and appearances. We use this data for both the control policy and vision-based state estimator. (b) The control policy receives observed robot states and rewards from the randomized simulations and learns to solve them using a recurrent neural network and reinforcement learning. (c) The vision-based state estimator uses rendered scenes collected from the randomized simulations and learns to predict the pose as well as face angles of the Rubik’s cube using a convolutional neural network (CNN), trained separately from the control policy. (d) To transfer to the real world, we predict the Rubik’s cube’s pose from 3 real camera feeds with the CNN and measure the robot fingertip locations using a 3D motion capture system. The face angles that describe the internal rotational state of the Rubik’s cube are provided by either the same vision state estimator *or* the Giiker cube, a custom cube with embedded sensors and feed it into the policy network.

has the potential to change this by *learning* how to use sensor information to control the robot system appropriately instead of hand-programming the robot using expert knowledge.

However, learning requires vast amount of training data, which is hard and expensive to acquire on a physical system. Collecting all data in simulation is therefore appealing. However, the simulation does not capture the environment or the robot accurately in every detail and therefore we also need to solve the resulting sim2real transfer problem. Domain randomization techniques [107, 81] have shown great potential and have demonstrated that models trained only in simulation can transfer to the real robot system.

In prior work, we have demonstrated that we can perform complex in-hand manipulation of a block [79]. This time, we aim to solve the manipulation and state estimation problems required to solve a Rubik’s cube with the Shadow Dexterous Hand [100] using only simulated data. This problem is much more difficult since it requires significantly more dexterity and precision for manipulating the Rubik’s cube. The state estimation problem is also much harder as we need to know with high accuracy what the pose and internal state of the Rubik’s cube are. We achieve this by introducing a novel method for automatically generating a distribution over randomized environments for training reinforcement learning policies and vision state estimators. We call this algorithm *automatic domain randomization* (ADR). We also built a robot platform for solving a Rubik’s cube in the real world in a way that complements our machine learning approach. Figure 2 shows an overview of our system.

We investigate why policies trained with ADR transfer so well from simulation to the real robot. We find clear signs of emergent learning that happens at *test time* within the recurrent internal state of our policy. We believe that this is a direct result of us training on an ever-growing distribution over randomized environments with a memory-augmented policy. In other words, training an LSTM over an ADR distribution is implicit meta-learning. We also systematically study and quantify this observation in our work.

The remainder of this manuscript is structured as follows. Section 2 introduces two manipulation tasks we consider here. Section 3 describes our physical setup and Section 4 describes how our setup is modeled in simulation. We introduce a new algorithm called automatic domain randomization (ADR), in Section 5. In Section 6 and Section 7 we describe how we train control policies and vision state estimators, respectively. We present our key quantitative and qualitative results on the two tasks in Section 8. In Section 9 we systematically analyze our policy for signs of emergent meta-learning. Section 10 reviews related work and we conclude with Section 11.

If you are mostly interested in the machine learning aspects of this manuscript, Section 5, Section 6, Section 7, Section 8, and Section 9 are especially relevant. If you are interested in the robotics aspects, Section 3, Section 4, and Section 8.4 are especially relevant.

## 2 Tasks

In this work, we consider two different tasks that both use the Shadow Dexterous Hand [100]: the block reorientation task from our previous work [79, 85] and the task of solving a Rubik’s cube. Both tasks are visualized in Figure 3. We briefly describe the details of each task in this section.

### 2.1 Block Reorientation

The block reorientation task was previously proposed in [85] and solved on a physical robot hand in [79]. We briefly review it here; please refer to the aforementioned citations for additional details.

The goal of the block reorientation task is to rotate a block into a desired goal orientation. For example, in Figure 3a, the desired orientation is shown next to the hand with the red face facing up, the blue face facing to the left and the green face facing forward. A goal is considered achieved if the block’s rotation matches the goal rotation within 0.4 radians. After a goal is achieved, a new random goal is generated.

### 2.2 Rubik’s Cube

We introduce a new and significantly more difficult problem in this work: solving a Rubik’s cube<sup>2</sup> with the same Shadow Dexterous Hand. In brief, a Rubik’s cube is a puzzle with 6 internal degrees of freedom. It consists of 26 *cubelets* that are connected via a system of joints and springs. Each of the 6 *faces* of the cube can be rotated, allowing the Rubik’s cube to be *scrambled*. A Rubik’s cube is considered solved if all 6 faces have been returned to a single color each. Figure 3b depicts a Rubik’s cube that is a single 90 degree rotation of the top face away from being solved.

<sup>2</sup>[https://en.wikipedia.org/wiki/Rubik's\\_Cube](https://en.wikipedia.org/wiki/Rubik's_Cube)

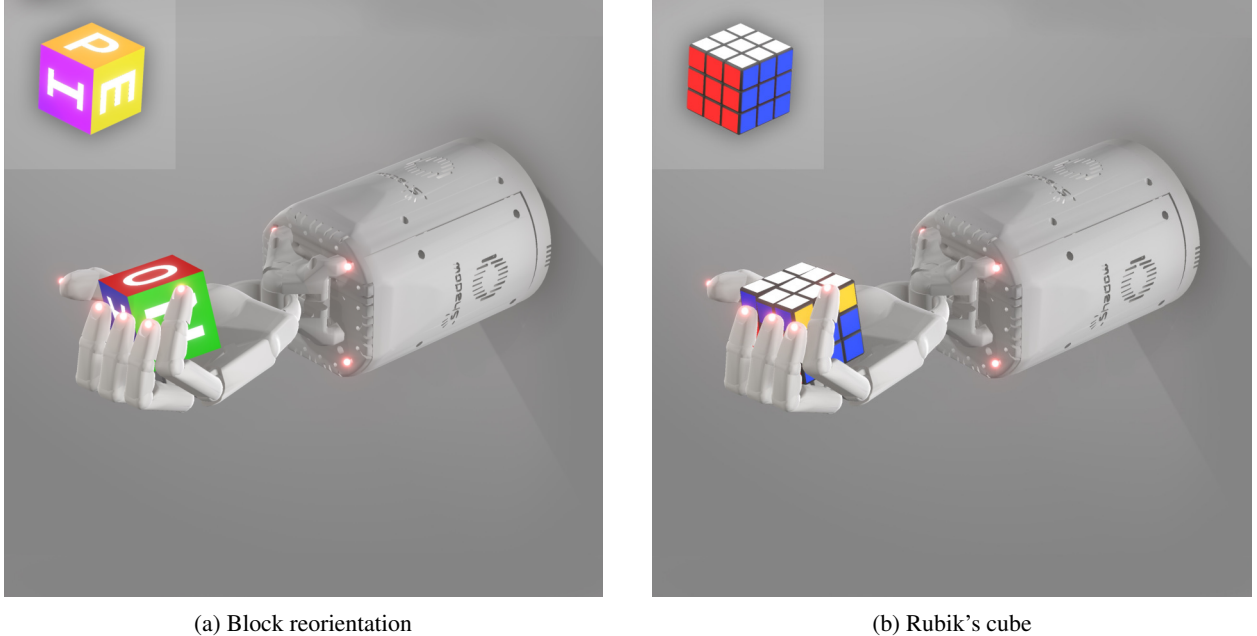


Figure 3: Visualization of the block reorientation task (left) and the Rubik's cube task (right). In both cases, we use a single Shadow Dexterous Hand to solve the task. We also depict the goal that the policy is asked to achieve in the upper left corner.

We consider two types of *subgoals*: A *rotation* corresponds to rotating a single face of the Rubik's cube by 90 degrees in the clockwise or counter-clockwise direction. A *flip* corresponds to moving a different face of the Rubik's cube to the top. We found rotating the top face to be far simpler than rotating other faces. Thus, instead of rotating arbitrary faces, we combine together a flip and a top face rotation in order to perform the desired operation. These subgoals can then be performed sequentially to eventually solve the Rubik's cube.

The difficulty of solving a Rubik's cube obviously depends on how much it has been scrambled before. We use the official scrambling method used by the World Cube Association<sup>3</sup> to obtain what they refer to as a *fair scramble*. A fair scramble typically consists of around 20 moves that are applied to a solved Rubik's cube to scramble it.

When it comes to solving the Rubik's cube, computing a solution sequence can easily be done with existing software libraries like the Kociemba solver [111]. We use this solver to produce a solution sequence of subgoals for the hand to perform. In this work, the key problem is thus about sensing and control, *not* finding the solution sequence. More concretely, we need to obtain the state of the Rubik's cube (i.e. its pose as well as its 6 face angles) and use that information to control the robot hand such that each subgoal is successfully achieved.

### 3 Physical Setup

Having described the task, we next describe the physical setup that we use to solve the block and the Rubik's cube in the real world. We focus on the differences that made it possible to solve the Rubik's cube since [79] has already described our physical setup for solving the block reorientation task.

#### 3.1 Robot Platform

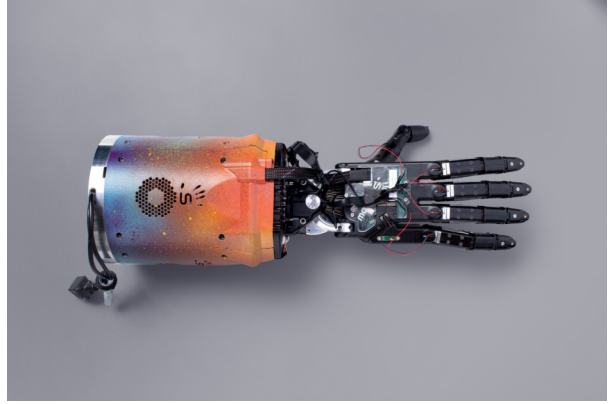
Our robot platform is based on the configuration described in [79]. We still use the Shadow Dexterous E Series Hand (E3M5R) [100] as a humanoid robot hand and the PhaseSpace motion capture system to track the Cartesian coordinates of all five fingertips. We use the same 3 RGB Basler cameras for vision pose estimation.

However, a number of improvements have been made since our previous publication. Figure 4a depicts the latest iteration of our robot cage. The cage is now fully contained, i.e. all computers are housed within the system. The cage is also on coasters and can therefore be moved more easily. The larger dimensions of the new cage make calibration of

<sup>3</sup><https://www.worldcubeassociation.org/regulations/scrambles/>



(a) The cage.



(b) The Shadow Dexterous Hand.

Figure 4: The latest version of our cage (left) that houses the Shadow Dexterous Hand, RGB cameras, and the PhaseSpace motion capture system. We made some modifications to the Shadow Dexterous Hand (right) to improve reliability for our setup by moving the PhaseSpace LEDs and cables inside the fingers and by adding rubber to the fingertips.

the PhaseSpace motion capture system easier and help prevent disturbing calibration when taking the hand in and out of the cage.

We have made a number of customizations to the E3M5R since our last publication (see also Figure 4b). We moved routing of the cables that connect the PhaseSpace LEDs on each fingertip to the PhaseSpace micro-driver within the hand, thus reducing the wear and tear on those cables. We worked with The Shadow Robot Company<sup>4</sup> to improve the robustness and reliability of some components for which we noticed breakages over time. We also modified the distal part of the fingers to extend the rubber area to cover a larger span to increase the grip of the hand when it interacts with an object. We increased the diameter of the wrist flexion/extension pulley in order to reduce tendon stress which has extended the life of the tendon to more than three times its typical mean time before failure (MTBF). Finally, the tendon tensioners in the hand have been upgraded and this has improved the MTBF of the finger tendons by approximately five to ten times.

We also made improvements to our software stack that interfaces with the E3M5R. For example, we found that manual tuning of the maximum torque that each motor can exercise was superior to our automated methods in avoiding physical breakage and ensuring consistent policy performance. More concretely, torque limits were minimized such that the hand can reliably achieve a series of commanded positions.

We also invested in real-time system monitoring so that issues with the physical setup could be identified and resolved more quickly. We describe our monitoring system in greater detail in Appendix A.

### 3.2 Giiker Cube

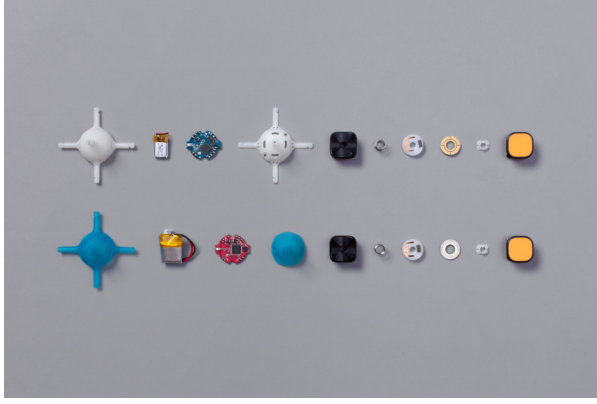
Sensing the state of a Rubik’s cube from vision only is a challenging task. We therefore use a “smart” Rubik’s cube with built-in sensors and a Bluetooth module as a stepping stone: We used this cube while face angle predictions from vision were not yet ready in order to continue work on the control policy. We also used the Giiker cube for some of our experiments to test the control policy without compounding errors made by the vision model’s face angle predictions (we always use the vision model for pose estimation).

Our hardware is based on the Xiaomi Giiker cube.<sup>5</sup> This cube is equipped with a Bluetooth module and allows us to sense the state of the Rubik’s cube. However, it only has a face angle resolution of  $90^\circ$ , which is not sufficient for state tracking purposes on the robot setup. We therefore replace some of the components of the original Giiker cube with custom ones in order to achieve a tracking accuracy of approximately  $5^\circ$ . Figure 5a shows the components of the unmodified Giiker cube and our custom replacements side by side, as well as the assembled modified Giiker cube. Since we only use our modified version, we henceforth refer to it as only “Giiker cube”.

<sup>4</sup><https://www.shadowrobot.com/>

<sup>5</sup><https://www.xiaomitoday.com/xiaomi-giiker-m3-intelligent-rubik-cube-review/>





(a) The components of the Giiker cube.



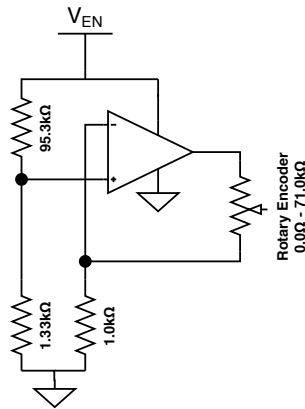
(b) An assembled Giiker cube while charging.

Figure 5: We use an off-the-shelf Giiker cube but modify its internals (subfigure a, right) to provide higher resolution for the 6 face angles. The components from left to right are (i) bottom center enclosure, (ii) lithium polymer battery, (iii) main PCBa with BLE, (iv) top center enclosure, (v) cubelet bottom, (vi) compression spring, (vii) contact brushes, (viii) absolute resistive rotary encoder, (ix) locking cap, (x) cubelet top. Once assembled, the Giiker cube can be charged with its “headphones on” (right).

### 3.2.1 Design

We have redesigned all parts of the Giiker cube but the exterior cubelet elements. The central support was redesigned to move the parting line off of the central line of symmetry to facilitate a more friendly development platform because the off-the-shelf design would have required de-soldering in order to program the microcontroller. The main Bluetooth and signal processing board is based on the NRF52 integrated circuit [75]. Six separately printed circuit boards (Figure 6b) were designed to improve the resolution from  $90^\circ$  to  $5^\circ$  using an absolute resistive encoder layout. The position is read with a linearizing circuit shown in Figure 6a. The linearized, analog signal is then read by an ADC pin on the microcontroller and sent as a face angle over the Bluetooth Low Energy (BLE) connection to the host.

The custom firmware implements a protocol that is based on the Nordic UART service (NUS) to emulate a serial port over BLE [75]. We then use a Node.js<sup>6</sup> based client application to periodically request angle readings from the UART module and to send calibration requests to reset angle references when needed. Starting from a solved Rubik’s cube, the client is able to track face rotations performed on the cube in real time and thus is able to reconstruct the Rubik’s cube state given periodic angle readings.



(a) The linearizing circuit used to read the position of the faces.



(b) The absolute resistive encoders used to read the position of the faces.

<sup>6</sup><https://nodejs.org/en/>

### 3.2.2 Data Accuracy and Timing

In order to ensure reliability of physical experiments, we performed regular accuracy tracking tests on integrated Giiker cubes. To assess accuracy, we considered all four right angle rotations as reference points on each cube face and estimated sensor accuracy based on measurements collected at each reference point. Across two custom cubes, the resistive encoders were subject to an absolute mean tracking error of  $5.90^\circ$  and the standard deviation of reference point readings was  $7.61^\circ$ .

During our experiments, we used a 12.5 Hz update frequency<sup>7</sup> for the angle readings, which was sufficient to provide low-latency observations to the robot policy.

### 3.2.3 Calibration

We perform a combination of firmware and software-side calibration of the sensors to ensure zero-positions can be dynamically set for each face angle sensor. On connecting to a cube for the first time, we record ADC offsets for each sensor in the firmware via a reset request. Furthermore, we add a software-side reset of the angle readings before starting each physical trial on the robot to ensure sensor errors do not accumulate across trials.

In order to track any physical degradation in the sensor accuracy of the fully custom hardware, we created a calibration procedure which instructs an operator to rotate each face a full  $360^\circ$ , stopping at each  $90^\circ$  alignment of the cube. We then record the expected and actual angles to measure the accuracy over time.

## 4 Simulation

The simulation setup is similar to [79]: we simulate the physical system with the MuJoCo physics engine [109], and we use ORRB [20], a remote rendering backend built on top of Unity3D,<sup>8</sup> to render synthetic images for training the vision based pose estimator.

While the simulation cannot perfectly match reality, we still found it beneficial to help bridge the gap by modeling our physical setup accurately. Our MuJoCo model of the Shadow Dexterous Hand has thus been further improved since [79] to better match the physical system via new dynamics calibration and modeling of a subset of tendons existing in the physical hand and we developed an accurate model of the Rubik’s cube.

### 4.1 Hand Dynamics Calibration

We measured joint positions for the same time series of actions for the real and simulated hands in an environment where the hand can move freely and made two observations:

1. The joint positions recorded on a physical robot and in simulation were visibly different (see Figure 8a).
2. The dynamics of *coupled joints* (i.e. distal two joints of non-thumb fingers, see [79, Appendix B.1]) were different on a physical robot and in simulation. In the original simulation used in [79], movement of coupled joints was modeled with two fixed tendons which resulted in both joints traveling roughly the same distance for each action. However, on the physical robot, movement of coupled joints depends on the current position of each joint. For instance, like in the human hand, the proximal segment of a finger bends before the distal segment when bending a finger.

To address the dynamics of coupled joints, we added a non-actuated spatial tendon and pulleys to the simulated non-thumb fingers (see Figure 7), analogous to the non-actuated tendon present in the physical robot. Parameters relevant to the joint movement in the new MuJoCo model were then calibrated to minimize root mean square error between reference joint positions recorded on a physical robot and joint positions recorded in simulation for the same time series of actions. We observe that better modeling of coupling and dynamics calibration improves performance significantly and present full results in Section D.1. We use this version of the simulation throughout the rest of this work.

### 4.2 Rubik’s Cube

Behind the apparent simplicity of a cube-like exterior, a Rubik’s cube hides a high degree of internal complexity and surprisingly nontrivial interactions between elements. A regular  $3 \times 3 \times 3$  cube consists of 26 externally facing *cubelets*

<sup>7</sup>We run the control policy at this frequency.

<sup>8</sup>Unity is a cross-platform game engine. See <https://www.unity.com> for more information.

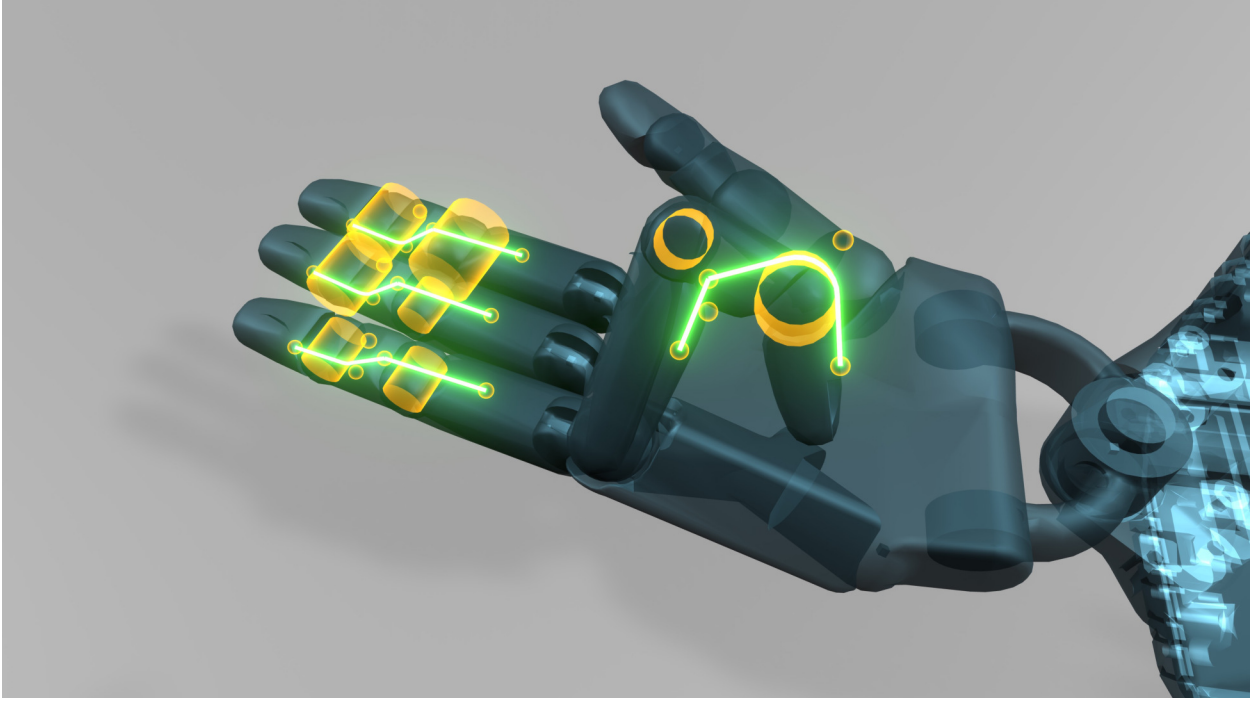


Figure 7: Transparent view of the hand in the new simulation. One spatial tendon (green lines) and two cylindrical geometries acting as pulleys (yellow cylinders) have been added for each non-thumb finger in order to achieve coupled joints dynamics similar to the physical robot.

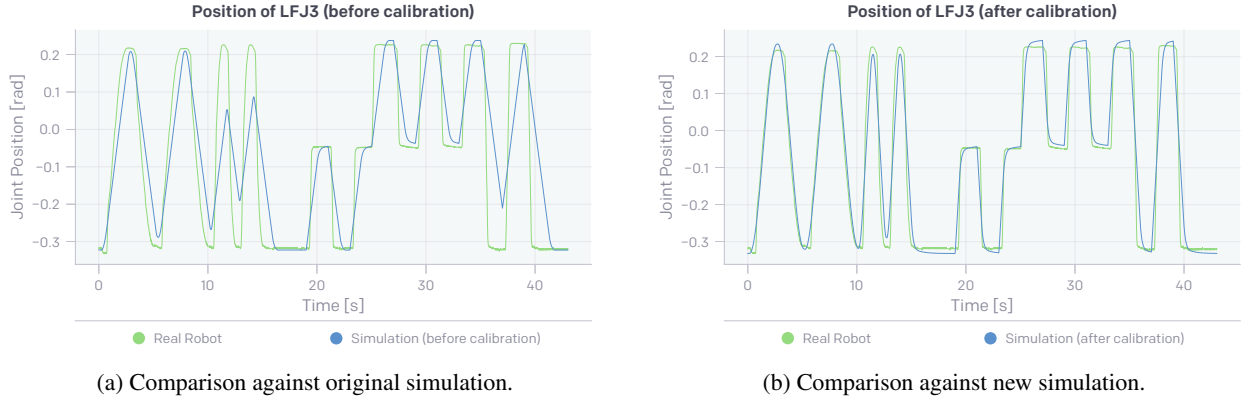
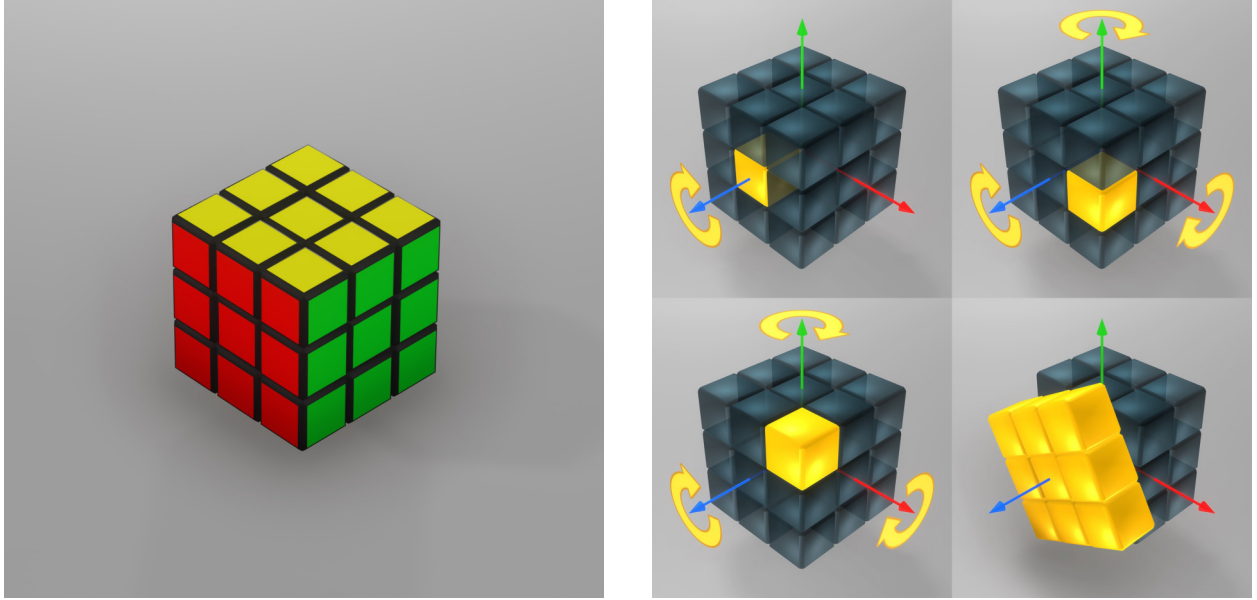


Figure 8: Comparison of positions of the LFJ3 joint on a real and simulated robot hand for the same control sequence, for the original simulation (a) and for the new simulation (b)

that are bound together to constitute a larger cubic shape. Six cubelets that reside in the center of each face are connected by axles to the inner core and can only rotate in place with one degree of freedom. In contrast to that, the edge and corner cubelets are not fixed and can move around the cube whenever the larger faces are rotated. To prevent the cube from falling apart, these cubelets have little plastic tabs that extend towards the core and allow each piece to be held in place by its neighbors, which in the end are retained by the center elements. Additionally, most Rubik's cubes are to a certain degree elastic and allow for small deformations from their original shape, constituting additional degrees of freedom.

The components of the cube constantly exert pressure on each other which results in a certain base level of friction in the system both between the cubelets and in the joints. It is enough to apply force to a single cubelet to rotate a face, as it will be propagated between the neighboring elements via contact forces. Although a cube has six faces that can be rotated, not all of them can be rotated simultaneously – whenever one face has already been moved by a certain





(a) Rendering of the cube model.

(b) Rendering of the different axis.

Figure 9: Our MuJoCo model of the Rubik's cube. On the left, we show a rendered version. On the right, we show the individual cubelets that make up our model and visualize the different axis and degrees of freedom of our model.

angle, perpendicular faces are in a locked state and prevented from moving. However, if this angle is small enough, the original face often "snaps" back into its nearest aligned state and in that way we can proceed with rotating the perpendicular face. This property is commonly called the "forgiveness" of a Rubik's Cube and its strength varies greatly among models available on the market.

Since we train entirely in simulation and need to successfully transfer to the real world without ever experiencing it, we needed to create a model rich enough to include all of the aforementioned behaviors, while at the same time keeping software complexity and computational costs manageable. We used the MuJoCo [109] physics engine, which implements a stable and fast numerical solutions for simulating body dynamics with soft contacts.

Inspired by the physical cube, our simulated model consists of 26 rigid body convex cubelets. MuJoCo allows for these shapes to penetrate each other by a small margin when a force is applied. Six central cubelets have a single *hinge joint* representing a single rotational degree of freedom about the axes running through the center of the cube orthogonal to each face. All remaining 20 corner and edge cubelets have three hinge joints corresponding to full Euler angle representation, with rotation axes passing through the center of the cube. In that way, our cube has  $6 \times 1 + 20 \times 3 = 66$  degrees of freedom, that allow us to represent effectively not only 43 quintillion fully aligned cube configurations but also all physically valid intermediate states.

Each cubelet mesh was created on the basis of the cube of size 1.9 cm. Our preliminary experiments have shown that with perfectly cubic shape, the overall Rubik's cube model was highly unforgiving. Therefore, we beveled all the edges of the mesh 1.425 mm inwards, which gave satisfactory results.<sup>9</sup> We do not implement any custom physics in our modelling, but rely on the cubelet shapes, contact forces and friction to drive the movement of the cube. We conducted experiments with spring joints which would correspond to additional degrees of freedom for cube deformation, but found they were not necessary and that native MuJoCo soft contacts already exhibit similar dynamics.

We performed a very rudimentary dynamics calibration of the parameters which MuJoCo allows us to specify, in order to roughly match a physical Rubik's cube. Our goal was not to get an exact match, but rather to have a plausible model as a starting point for domain randomization.

<sup>9</sup>Real Rubik's cubes also have cubelets with rounded corners, for the same reason.

## 5 Automatic Domain Randomization

In [79], we were able to train a control policy and a vision model in simulation and then transfer both to a real robot through the use of domain randomization [107, 81]. However, this required a significant amount of manual tuning and a tight iteration loop between randomization design in simulation and validation on a robot. In this section, we describe how *automatic domain randomization* (ADR) can be used to automate this process and how we apply ADR to both policy and vision training.

Our main hypothesis that motivates ADR is that *training on a maximally diverse distribution over environments leads to transfer via emergent meta-learning*. More concretely, if the model has some form of memory, it can learn to adjust its behavior during deployment to improve performance on the current environment over time, i.e. by implementing a learning algorithm internally. We hypothesize that this happens if the training distribution is so large that the model cannot memorize a special-purpose solution per environment due to its finite capacity. ADR is a first step in this direction of unbounded environmental complexity: it automates and gradually expands the randomization ranges that parameterize a distribution over environments. Related ideas were also discussed in [31, 15, 117, 24].

In the remainder of this section, we first describe how ADR works at a high level and then describe the algorithm and our implementation in greater detail.

### 5.1 ADR Overview

We use ADR both to train our vision models (supervised learning) and our policy (reinforcement learning). In each case, we generate a distribution over environments by randomizing certain aspects, e.g. the visual appearance of the cube or the dynamics of the robotic hand. While domain randomization requires us to define the ranges of this distribution manually and keep it fixed throughout model training, in ADR the distribution ranges are defined automatically and allowed to change.

A top-level diagram of ADR is given in Figure 10. We give an intuitive overview of ADR below. See Section 5.2 for a formal description of the algorithm.

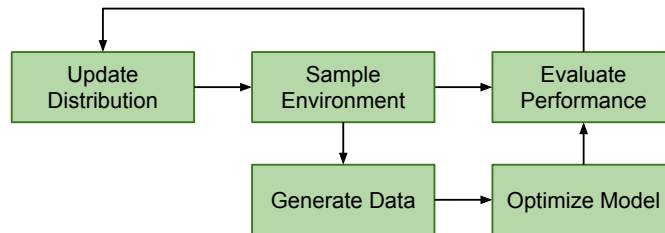


Figure 10: Overview of ADR. ADR controls the distribution over environments. We sample environments from this distribution and use it to generate training data, which is then used to optimize our model (either a policy or a vision state estimator). We further evaluate performance of our model on the current distribution and use this information to update the distribution over environments automatically.

At its core, ADR realizes a training curriculum that gradually expands a distribution over environments for which the model can perform well. The initial distribution over environments is concentrated on a single environment. For example, in policy training the initial environment is based on calibration values measured from the physical robot.

The distribution over environments is sampled to obtain environments used to generate training data and evaluate model performance. ADR is independent of the algorithm used for model training. It only generates training data. This allows us to use ADR for both policy and vision model training.

As training progresses and model performance improves sufficiently on the initial environment, the distribution is expanded. This expansion continues as long as model performance is considered acceptable. With a sufficiently powerful model architecture and training algorithm, the distribution is expected to expand far beyond manual domain randomization ranges since every improvement in the model’s performance results in an increase in randomization.

ADR has two key benefits over manual domain randomization (DR):

- Using a curriculum that gradually increases difficulty as training progresses simplifies training, since the problem is first solved on a single environment and additional environments are only added when a minimum level of performance is achieved [37, 69].
- It removes the need to manually tune the randomizations. This is critical, because as more randomization parameters are incorporated, manual adjustment becomes increasingly difficult and non-intuitive.

Acceptable performance is defined by *performance thresholds*. For policy training, they are configured as the lower and upper bounds on the number of successes in an episode. For vision training, we first configure target performance thresholds for each output (e.g. position, orientation). During evaluation, we then compute the percentage of samples which achieve these targets for all outputs; if the resulting percentage is above the upper threshold or below the lower threshold, the distribution is adjusted accordingly.

## 5.2 Algorithm

Each environment  $e_\lambda$  is parameterized by  $\lambda \in \mathbb{R}^d$ , where  $d$  is the number of parameters we can randomize in simulation. In domain randomization (DR), the environment parameter  $\lambda$  comes from a *fixed* distribution  $P_\phi$  parameterized by  $\phi \in \mathbb{R}^{d'}$ . However, in automatic domain randomization (ADR),  $\phi$  is *changing* dynamically with training progress. The sampling process in Figure 10 works out as  $\lambda \sim P_\phi$ , resulting in one randomized environment instance  $e_\lambda$ .

To quantify the amount of ADR expansion, we define *ADR entropy* as  $\mathcal{H}(P_\phi) = -\frac{1}{d} \int P_\phi(\lambda) \log P_\phi(\lambda) d\lambda$  in units of nats/dimension. The higher the ADR entropy, the broader the randomization sampling distribution. The normalization allows us to compare between different environment parameterizations.

In this work, we use a factorized distribution parameterized by  $d' = 2d$  parameters. To simplify notation, let  $\phi^L, \phi^H \in \mathbb{R}^d$  be a certain partition of  $\phi$ . For the  $i$ -th ADR parameter  $\lambda_i$ ,  $i = 1, \dots, d$ , the pair  $(\phi_i^L, \phi_i^H)$  is used to describe a uniform distribution for sampling  $\lambda_i$  such that  $\lambda_i \sim U(\phi_i^L, \phi_i^H)$ . Note that the boundary values are inclusive. The overall distribution is given by

$$P_\phi(\lambda) = \prod_{i=1}^d U(\phi_i^L, \phi_i^H)$$

with ADR entropy

$$\mathcal{H}(P_\phi) = \frac{1}{d} \sum_{i=1}^d \log(\phi_i^H - \phi_i^L).$$

The ADR algorithm is listed in Algorithm 1. For the factorized distribution, Algorithm 1 is applied to  $\phi^L$  and  $\phi^H$  separately.

At each iteration, the ADR algorithm randomly selects a dimension of the environment  $\lambda_i$  to fix to a boundary value  $\phi_i^L$  or  $\phi_i^H$  (we call this “boundary sampling”), while the other parameters are sampled as per  $P_\phi$ . Model performance for the sampled environment is then evaluated and appended to the buffer associated with the selected boundary of the selected parameter. Once enough performance data is collected, it is averaged and compared to thresholds. If the average model performance is better than the high threshold  $t_H$ , the parameter for the chosen dimension is increased. It is decreased if the average model performance is worse than the low threshold  $t_L$ .

As described, the ADR algorithm modifies  $P_\phi$  by always fixing one environment parameter to a boundary value. To generate model training data, we use Algorithm 2 in conjunction with ADR. The algorithm samples  $\lambda$  from  $P_\phi$  and runs the model in the sampled environment to generate training data.

To combine ADR and training data generation, at every iteration we execute Algorithm 1 with probability  $p_b$  and Algorithm 2 with probability  $1 - p_b$ . We refer to  $p_b$  as the *boundary sampling probability*.

## 5.3 Distributed Implementation

We used a distributed version of ADR in this work. The system architecture is illustrated in Figure 11 for both our policy and vision training setup. We describe policy training in greater detail in Section 6 and vision training in Section 7. Here we focus on ADR.

The highly-parallel and asynchronous implementation depends on several centralized storage of (policy or vision) model parameters  $\Theta$ , ADR parameters  $\Phi$ , training data  $T$ , and performance data buffers  $\{D_i\}_{i=1}^d$ . We use Redis to implement them.

**Algorithm 1** ADR

---

**Require:**  $\phi^0$  ▷ Initial parameter values  
**Require:**  $\{D_i^L, D_i^H\}_{i=1}^d$  ▷ Performance data buffers  
**Require:**  $m, t_L, t_H$ , where  $t_L < t_H$  ▷ Thresholds  
**Require:**  $\Delta$  ▷ Update step size  
 $\phi \leftarrow \phi^0$   
**repeat**  
     $\lambda \sim P_\phi$   
     $i \sim U\{1, \dots, d\}, x \sim U(0, 1)$   
    **if**  $x < 0.5$  **then**  
         $D_i \leftarrow D_i^L, \lambda_i \leftarrow \phi_i^L$  ▷ Select the lower bound in “boundary sampling”  
    **else**  
         $D_i \leftarrow D_i^H, \lambda_i \leftarrow \phi_i^H$  ▷ Select the higher bound in “boundary sampling”  
    **end if**  
     $p \leftarrow \text{EVALUATEPERFORMANCE}(\lambda)$  ▷ Collect model performance on environment parameterized by  $\lambda$   
     $D_i \leftarrow D_i \cup \{p\}$  ▷ Add performance to buffer for  $\lambda_i$ , which was boundary sampled  
    **if**  $\text{LENGTH}(D_i) \geq m$  **then**  
         $\bar{p} \leftarrow \text{AVERAGE}(D_i)$   
         $\text{CLEAR}(D_i)$   
        **if**  $\bar{p} \geq t_H$  **then**  
             $\phi_i \leftarrow \phi_i + \Delta$   
        **else if**  $\bar{p} \leq t_L$  **then**  
             $\phi_i \leftarrow \phi_i - \Delta$   
        **end if**  
    **end if**  
**until** training is complete

---

**Algorithm 2** Training Data Generation

---

**Require:**  $\phi$  ▷ ADR distribution parameters  
**repeat**  
     $\lambda \sim P_\phi$   
     $\text{GENERATEDATA}(\lambda)$   
**until** training is complete

---

By using centralized storage, the ADR algorithm is decoupled from model optimization. However, to train a good policy or vision model using ADR, it is necessary to have a concurrent optimizer that consumes the training data in  $T$  and pushes updated model parameters to  $\Theta$ .

We use  $W$  parallel worker threads instead of the sequential while-loop. For training the policy, each worker pulls the latest distribution and model parameters from  $\Phi$  and  $\Theta$  and executes Algorithm 1 with probability  $p_b$  (denoted as “ADR Eval Worker” in Figure 11a). Otherwise, it executes Algorithm 2 and pushes the generated data to  $T$  (denoted as “Rollout Worker” in Figure 11a). To avoid wasting a large amount of data for only ADR, we also use this data to train the policy. The setup for vision is similar. Instead of rolling out a policy, we use the ADR parameters to render images and use those to train the supervised vision state estimator. Since data is cheaper to generate, we do not use the ADR evaluator data to train the model in this case but only used the data produced by the “Data Producer” (compare Figure 11b).

In the policy model,  $\phi^0$  is set based on a calibrated environment parameter according to  $\phi_i^{0,L} = \phi_i^{0,H} = \lambda_i^{\text{calib}}$  for all  $i = 1, \dots, d$ . In the vision model, the initial randomizations are set to zero, i.e.  $\phi_i^{0,L} = \phi_i^{0,H} = 0$ . The distribution parameters are pushed to  $\Phi$  to be used by all workers at the beginning of the algorithm.

## 5.4 Randomizations

Here, we describe the categories of randomizations used in this work. The vast majority of randomizations are for a scalar environment parameter  $\lambda_i$  and are parameterized in ADR by two boundary parameters  $(\phi_i^L, \phi_i^H)$ . For a full listing of randomizations used in policy and vision training, see Appendix B.

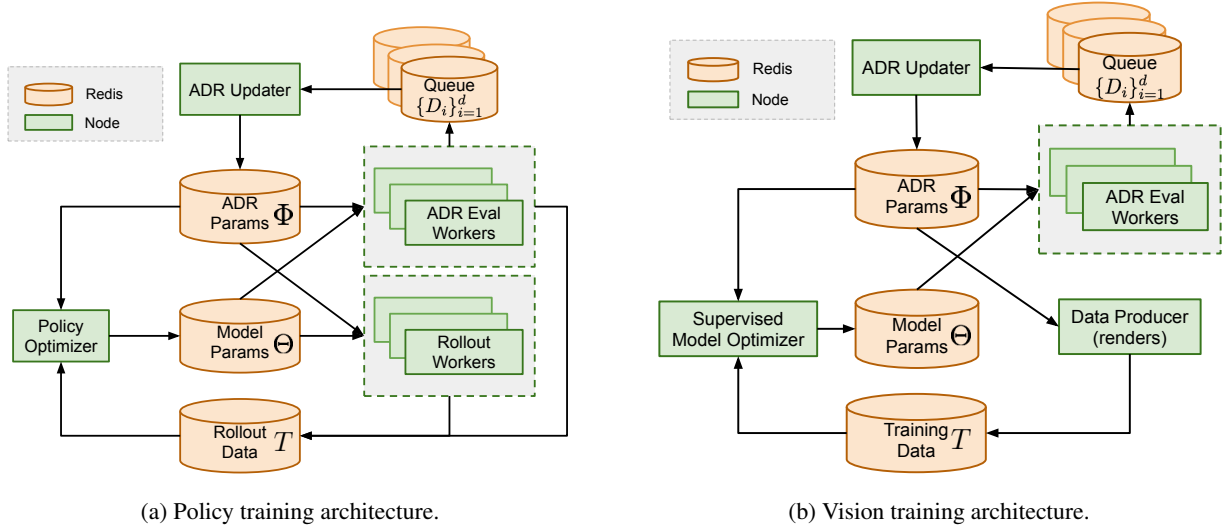


Figure 11: The distributed ADR architecture for policy (left) and vision (right). In both cases, we use Redis for centralized storage of ADR parameters ( $\Phi$ ), model parameters ( $\Theta$ ), and training data ( $T$ ). ADR eval workers run Algorithm 1 to estimate performance using boundary sampling and report results using performance buffers ( $\{D_i\}_{i=1}^d$ ). The ADR updater uses those buffers to obtain average performance and increases or decreases boundaries accordingly. Rollout workers (for the policy) and data producers (for vision) produce data by sampling an environment as parameterized by the current set of ADR parameters (see Algorithm 2). This data is then used by the optimizer to improve the policy and vision model, respectively.

A few randomizations, such as observation noise, are controlled by more than one environment parameter and are parameterized by a larger set of boundary parameters. For full details on these randomizations and their ADR parameterization, see Appendix B.

**Simulator physics.** We randomize simulator physics parameters such as geometry, friction, gravity, etc. See Section B.1 for details of their ADR parameterization.

**Custom physics.** We model additional physical robot effects that are not modelled by the simulator, for example, action latency or motor backlash. See [79, Appendix C.2] for implementation details of these models. We randomize the parameters in these models in a similar way to simulator physics randomizations.

**Adversarial.** We use an adversarial approach similar to [83, 84] to capture any remaining unmodeled physical effects in the target domain. However, we use random networks instead of a trained adversary. See Section B.3 for details on implementation and ADR parameterization.

**Observation.** We add Gaussian noise to policy observations to better approximate observation conditions in reality. We apply both correlated noise, which is sampled once at the start of an episode and uncorrelated noise, which is sampled at each time step. We randomize the parameters of the added noise. See Section B.4 for details of their ADR parameterization.

**Vision.** We randomize several aspects in ORRB [20] to control the rendered scene, including lighting conditions, camera positions and angles, materials and appearances of all the objects, the texture of the background, and the post-processing effects on the rendered images. See Section B.5 for details.

## 6 Policy Training in Simulation

In this section we describe how we train control policies using Proximal Policy Optimization [99] and reinforcement learning. Our setup is similar to [79]. However, we use ADR as described in Section 5 to train on a large distribution over randomized environments.



## 6.1 Actions, Rewards, and Goals

Our setup for the action space and rewards is unchanged from [79] so we only briefly recap them here. We use a discretized action space with 11 bins per actuated joint (of which there are 20). We use a multi-categorical distribution. Actions are relative changes in generalized joint position coordinates.

There are three types of rewards we provide to our agent during training: (a) The difference between the previous and the current distance of the system state from the goal state, (b) an additional reward of 5 whenever a goal is achieved, (c) and a penalty of  $-20$  whenever a cube/block is dropped.

We generate random goals during training. For the block, the target rotation is randomly sampled but constrained such that any face points directly upwards. For the Rubik’s cube the task generation is slightly more convoluted as it depends on the state of the cube at the time when the goal is generated. If the cube faces are not aligned, we make sure to align them and additionally rotate the whole cube according to a sampled random orientation just like with the block (called a flip). Alternatively, if the faces *are* aligned, we rotate the top cube face with 50% probability either clockwise or counter-clockwise. Otherwise we again perform a flip. Detailed listings of the goal generation algorithms can be found in the Section C.1.

We consider a training episode to be finished whenever one of the following conditions is satisfied: (a) the agent achieves 50 consecutive successes (of reaching a goal within the required threshold), (b) the agent drops the cube, (c) or the agent times out when trying to reach the next goal. Time out limits are 400 timesteps for block reorientation and 800 timesteps<sup>10</sup> for the Rubik’s Cube.

## 6.2 Policy Architecture

We base our policy architecture on [79] but extend it in a few important ways. The policy is still recurrent since only a policy with access to some form of memory can perform meta-learning. We still use a single feed-forward layer with a ReLU activation [74] followed by a single LSTM layer [47]. However, we increase the capacity of the network by doubling the number of units: the feed-forward layer now has 2048 units and the LSTM layer has 1024 units.

The value network is separate from the policy network (but uses the same architecture) and we project the output of the LSTM onto a scalar value. We also add L2 regularization with a coefficient of  $10^{-6}$  to avoid ever-growing weight norms for long-running experiments.

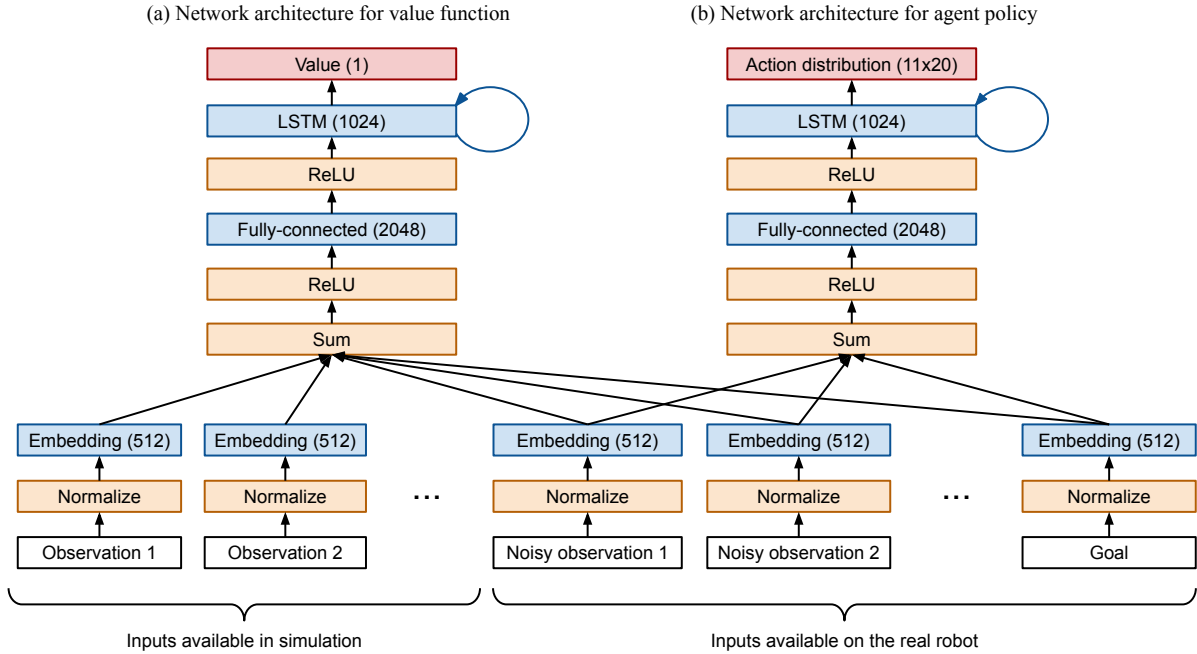


Figure 12: Neural network architecture for (a) value network and (b) policy network.

<sup>10</sup>We use 1600 timesteps when training from scratch.

Table 1: Inputs for the Rubik’s cube task of the policy and value networks, respectively.

Input	Dimensionality	Policy network	Value network
Fingertip positions	15D	×	✓
Noisy fingertip positions	15D	✓	✓
Cube position	3D	×	✓
Noisy cube position	3D	✓	✓
Cube orientation	4D (quaternion)	×	✓
Noisy cube orientation	4D (quaternion)	✓	✓
Goal orientation	4D (quaternion)	✓	✓
Relative goal orientation	4D (quaternion)	×	✓
Noisy relative goal orientation	4D (quaternion)	✓	✓
Goal face angles	12D <sup>11</sup>	✓	✓
Relative goal face angles	12D <sup>11</sup>	×	✓
Noisy relative goal face angles	12D <sup>11</sup>	✓	✓
Hand joint angles	48D <sup>11</sup>	×	✓
All simulation positions & orientations (qpos)	170D	×	✓
All simulation velocities (qvel)	168D	×	✓

An important difference between our architecture and the architecture used in [79] is how inputs are handled. In [79], the inputs for the policy and value networks consisted of different observations (e.g. fingertip positions, block pose, ...) in noisy and non-noisy versions. For each network, all observation fields were concatenated into a single vector. Noisy observations were provided to the policy network while the value network had access to non-noisy observations (since the value network is not needed when rolling out the policy on the robot and can thus use privileged information, as described in [82]). We still use the same Asymmetric Actor-Critic architecture [82] but replace the concatenation with what we call an “embed-and-add” approach. More concretely, we first embed each type of observation separately (without any weight sharing) into a latent space of dimensionality 512. We then combine all inputs by adding the latent representation of each and applying a ReLU non-linearity after. The main motivation behind this change was to easily add new observations to an existing policy and to share embeddings between value and policy network for inputs that feed into both. The network architecture of our control policy is illustrated in Figure 12. More details of what inputs are fed into the networks can be found in Section C.3 and Table 1. We list the inputs for the block reorientation task in Section C.2.

### 6.3 Distributed Training with Rapid

We use our own internal distributed training framework, Rapid. Rapid was previously used to train OpenAI Five [78] and was also used in [79].

For the block reorientation task, we use  $4 \times 8 = 32$  NVIDIA V100 GPUs and  $4 \times 100 = 400$  worker machines with 32 CPU cores each. For the Rubik’s cube task, we use  $8 \times 8 = 64$  NVIDIA V100 GPUs and  $8 \times 115 = 920$  worker machines with 32 CPU cores each. We’ve been training the Rubik’s Cube policy continuously for several months at this scale while concurrently improving the simulation fidelity, ADR algorithm, tuning hyperparameters, and even changing the network architecture. The cumulative amount of experience over that period used for training on the Rubik’s cube is roughly 13 thousand years, which is on the same order of magnitude as the 40 thousand years used by OpenAI Five [78].

The hyperparameters that we used and more details on optimization can be found in Section C.3.

### 6.4 Policy Cloning

With ADR, we found that training the same policy for a very long time is helpful since ADR allows us to always have a challenging training distribution. We therefore rarely trained experiments from scratch but instead updated existing experiments and initialized from previous checkpoints for both the ADR and policy parameters. Our new “embed-and-add” approach in Figure 12 makes it easier to change the observation space of the agent, but doesn’t allow

<sup>11</sup>Angles are encoded as sin and cos, i.e. this doubles the dimensionality of the underlying angle.

us to experiment with changes to the policy architecture, e.g. modify the number of units in each layer or add a second LSTM layer. Restarting training from an uninitialized model would have caused us to lose weeks or months of training progress, making such changes prohibitively expensive. Therefore, we successfully implemented behavioral cloning in the spirit of the DAGGER [89] algorithm (sometimes also called policy distillation [27]) to efficiently initialize new policies with a level of performance very close to the teacher policy.

Our setup for cloning closely mimics reinforcement learning, except that we now have both teacher and student policies loaded in memory. During a rollout, we use the student actions to interact with the environment, while minimizing the difference between the student and the teacher’s action distributions (by minimizing KL divergence) and value predictions (by minimizing L2 loss). This has worked surprisingly well, allowing us to iterate on the policy architecture quickly without losing the accumulated training progress. Our cloning approach works with arbitrary policy architecture changes as long as the action space remains unchanged.

The best ADR policies used in this work were obtained using this approach. We trained them for multiple months while making multiple changes to the model architecture, training environment, and hyperparameters.

## 7 State Estimation from Vision

As in [79], the control policy described in Section 6 receives object state estimates from a vision system consisting of three cameras and a neural network predictor. In this work, the policy requires estimates for all six face angles in addition to the position and orientation of the cube.

Note that the absolute rotation of each face angle in  $[-\pi, \pi]$  radians is required by the policy. Due to the rotational symmetry of the stickers on a standard Rubik’s cube, it is not possible to predict these absolute face angles from a single camera frame; the system must either have some ability to track state temporally<sup>12</sup> or the cube has to be modified.

We therefore use two different options for the state estimation of the Rubik’s cube throughout this work:

1. **Vision only via asymmetric center stickers.** In this case, the vision model is used to produce the *cube position, rotation, and six face angles*. We cut out one corner of each center sticker on the cube (see Figure 13), thus breaking rotational symmetry and allowing our model to determine absolute face angles from a single frame. No further customizations were made to the Rubik’s cube. We use this model to estimate final performance of a vision only solution to solving the Rubik’s cube.
2. **Vision for pose and Giiker cube for face angles.** In this case, the vision model is used to produce the *cube position and rotation*. For the face angles, we use the previously described customized Giiker cube (see Section 3) with built-in sensors. We use this model for most experiments in order to not compound errors of the challenging face angle estimation from vision only with errors of the policy.

Since our long-term goal is to build robots that can interact in the real world with arbitrary objects, ideally we would like to fully solve this problem from vision alone using a standard Rubik’s cube (i.e. without any special stickers). We believe this is possible, though it may require either more extensive work on a recurrent model or moving to an end-to-end training setup (i.e. where the vision model is learned jointly with the policy). This remains an active area of research for us.

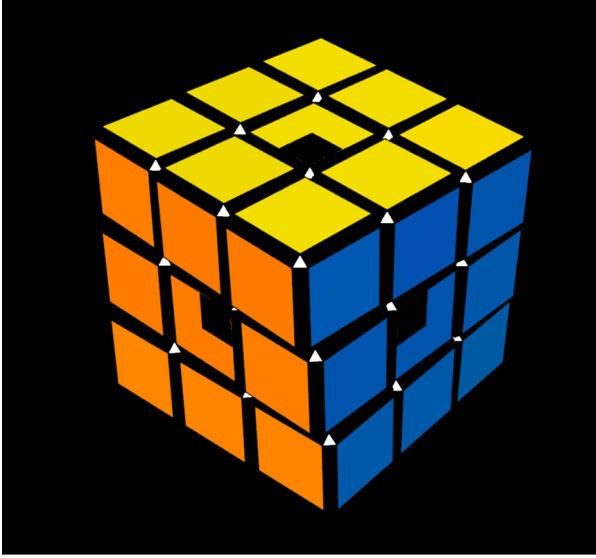
### 7.1 Vision Model

Our vision model has a similar setup as in [79], taking as input an image from each of three RGB Basler cameras located at the left, right, and top of the cage (see Figure 4(a)). The full model architecture is illustrated in Figure 14. We produce a feature map for each image by processing it through identically parameterized ResNet50 [45] networks (i.e. using common weights). These three feature maps are then flattened, concatenated, and fed into a stack of fully-connected layers which ultimately produce predictions sufficient for tracking the full state of the cube, including the position, orientation, and face angles.

While predicting position and orientation directly works well, we found predicting all six face angles directly to be much more challenging due to heavy occlusion, even when using a cube with asymmetric center stickers. To work around this, we decomposed face angle prediction into several distinct predictions:

1. **Active axis:** We make a slight simplifying assumption that only one of the three axes of a cube can be "active" (i.e. be in a non-aligned state), and have the model predict which of the three axes is currently active.

<sup>12</sup>We experimented with a recurrent vision model but found it very difficult to train to the necessary performance level. Due to the project’s time constraints, we could not investigate this approach further.



(a) Simulated cube.



(b) Real cube.

Figure 13: The Rubik’s cube with a corner cut out of each center sticker (a) in simulation and (b) in reality. We used this cube instead of the Giiker cube for some vision state estimation experiments and for evaluating the performance of the policy for solving the Rubik’s cube from vision only.

2. **Active face angles:** We predict the angles of the two faces relevant for the active axis *modulo*  $\pi/2$  radians (i.e. in  $[-\pi/4, \pi/4]$ ). It is hard to predict the absolute angles in  $[-\pi, \pi]$  radians directly due to heavy occlusion (e.g. when a face is on the bottom and hidden by the palm). Predicting these modulo  $\pi/2$  angles only requires recognizing the shape and the relative positions of cube edges, and therefore it is an easier task.
3. **Top face angle:** The last piece to predict is the absolute angle in  $[-\pi, \pi]$  radians of the "top" face, that is the face visible from a camera mounted directly above the hand. Note that this angle is only possible to predict from frames at a single timestamp because of the asymmetric center stickers (See Figure 13). We configure the model to make a prediction only for the top face because the top face’s center cubelet is rarely occluded. This gives us a stateless estimate of each face’s absolute angle of rotation whenever that face is placed on top.

These decomposed face angle predictions are then fed into post-processing logic (See Appendix C Algorithm 5) to track the rotation of all face angles, which are in turn passed along to the policy. The top face angle prediction is especially important, as it allows us to correct the tracked absolute face angle state mid-trial. For example, if the tracking of a face angle becomes off by some number of rotations (i.e. a multiple of  $\pi/2$  radians), we are still able to correct it with a stateless absolute angle prediction from the model whenever this face is placed on top after a flip. Predictions (1) and (2) are primarily important because the policy is unable to rotate a non-active face if the active face angles are too large (in which case the cube becomes interlocked along non-active axes).

For all angle predictions, we found that discretizing angles into 90 bins per  $\pi$  radians yielded better performance than directly predicting angles via regression; see Table 2 for details.

In the meantime, domain randomization in the rendering process remains a critical role in the sim2real transfer. As shown in Table 2, a model trained without domain randomization can achieve perfectly low errors in simulation but fails dramatically on real world data.

## 7.2 Distributed Training with Rapid

As in control policy training (Section 6), the vision model is trained entirely from synthetic data, without any images from the real world. This necessarily entails a more complicated training setup, wherein the synthetic image generation must be coupled with optimization. To manage this complexity, we leverage the same Rapid framework [78] which is used in policy training for distributed training.

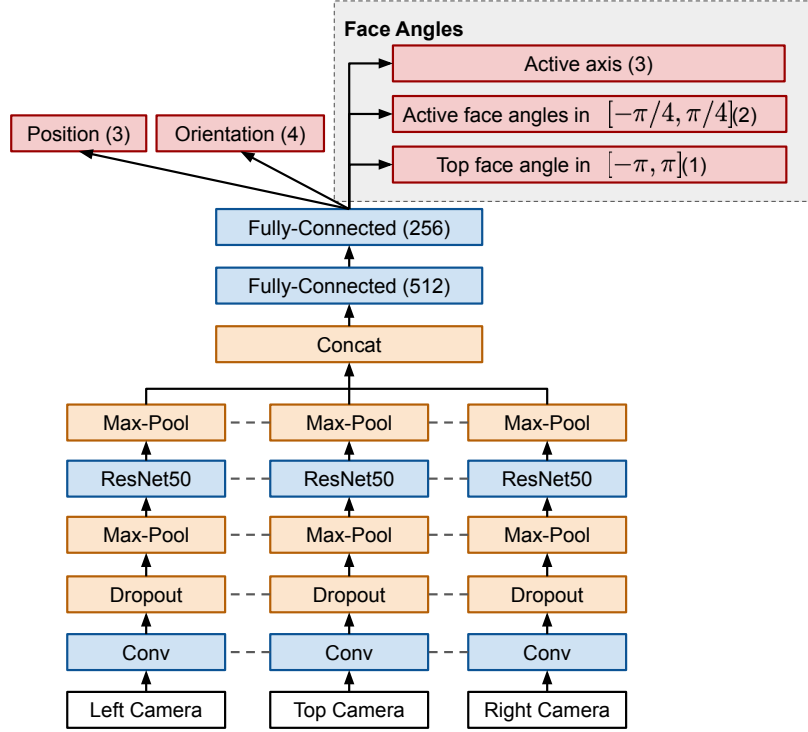


Figure 14: Vision model architecture, which is largely built upon a ResNet50 [45] backbone. Network weights are shared across the three camera frames, as indicated by the dashed line. Our model produces the position, orientation, and a specific representation of the six face angles of the Rubik’s cube. We specify ranges with  $[\dots]$  and dimensionality with  $(\dots)$ .

Table 2: Ablation experiments for the vision model. For each experiment, we ran training with 3 different seeds and report the best performance here. Orientation error is computed as rotational distance over a quaternion representation. Position error is the euclidean distance in 3D space, in millimeters. Face angle error is measured in degrees ( $^\circ$ ). "Real" errors are computed using data collected over multiple physical trials, where the position and orientation ground truths are from PhaseSpace (Section 3) and all face angle ground truths are from the Giiker cube. The full evaluation results, including errors on active axis and active face angles, are reported in Appendix D Table 22.

Experiment	Errors (Sim)			Errors (Real)		
	Orientation	Position	Top Face	Orientation	Position	Top face
Full Model	6.52 $^\circ$	<b>2.63 mm</b>	11.95 $^\circ$	<b>7.81<math>^\circ</math></b>	<b>6.47 mm</b>	<b>15.92<math>^\circ</math></b>
No Domain Randomization	<b>3.95<math>^\circ</math></b>	2.97 mm	<b>8.56<math>^\circ</math></b>	128.83 $^\circ$	69.40 mm	85.33 $^\circ$
No Focal Loss	15.94 $^\circ$	5.02 mm	10.17 $^\circ$	19.10 $^\circ$	9.416 mm	17.54 $^\circ$
Non-discrete Angles	9.02 $^\circ$	3.78 mm	42.46 $^\circ$	10.40 $^\circ$	7.97 mm	35.27 $^\circ$

Figure 11b gives an overview of the setup for a typical vision experiment. In the case of vision training, the “data workers” are standalone Unity renderers, responsible for rendering simulated images using OpenAI Remote Rendering Backend (ORRB) [20]. These images are rendered according to ADR parameters pulled from the ADR subsystem (see Section 5). A list of randomization parameters is available in Section B.5 Table 11. Each rendering node uses 1 NVIDIA V100 GPU and 8 CPU cores, and the size of the rendering pool is tuned such that rendering is not a bottleneck in training. The data from these rendering nodes is then propagated to a cluster of Redis nodes where it is stored in separate queues for training and evaluation. The training data is then read by a pool of optimizer nodes, each of which



uses 8 NVIDIA V100 GPUs and 64 CPU cores, in order to perform optimization in a data-parallel fashion. Meanwhile, the evaluation data is read by the “ADR eval workers” in order to provide feedback on ADR parameters, per Section 5.

As noted above, the vision model produces several distinct predictions, each of which has its own loss function to be optimized: mean squared error for both position and orientation, and cross entropy for each of the decomposed face angle predictions. To balance these many losses, which lie on different scales, we use focal loss weighting as described in [39] to dynamically and automatically assign loss weights. One modification we made in order to better fit our multiple regression tasks is that we define a low target error for each prediction and then use the percentage of samples that obtain errors below the target as the probability  $p$  in the focal loss, i.e.  $FL(p; \gamma) = -(1 - p)^\gamma \log(p)$ , where  $\gamma = 1$  in all our experiments. This both removes the need to manually tune loss weights and improves optimization, as it allows loss weights to change dynamically during training (see Table 2 for performance details).

Optimization is then performed against this aggregate loss using the LARS optimizer [118]. We found LARS to be more stable than the Adam optimizer [53] when using larger batches and higher learning rates (we use at most a batch size of 1024 with a peak learning rate of 6.0). See Section C.3 for further hyperparameter details.

## 8 Results

In this section, we investigate the effect ADR has on transfer (Section 8.1), empirically show the importance of having a curriculum for policy training (Section 8.2), quantify vision performance (Section 8.3), and finally present our results that push the limits of what is possible by solving a Rubik’s cube on the real Shadow hand (Section 8.4).

### 8.1 Effect of ADR on Policy Transfer

To understand the transfer performance impact of training policies with ADR, we study the problem on the simpler block reorientation task previously introduced in [79]. We use this task since it is computationally more tractable and because baseline performance has been established. As in [79], we measure performance in terms of the number of consecutive successes. We terminate an episode if the block is either dropped or if 50 consecutive successes are achieved. An optimal policy would therefore be one that achieves a mean of 50 successes.

#### 8.1.1 Sim2Sim

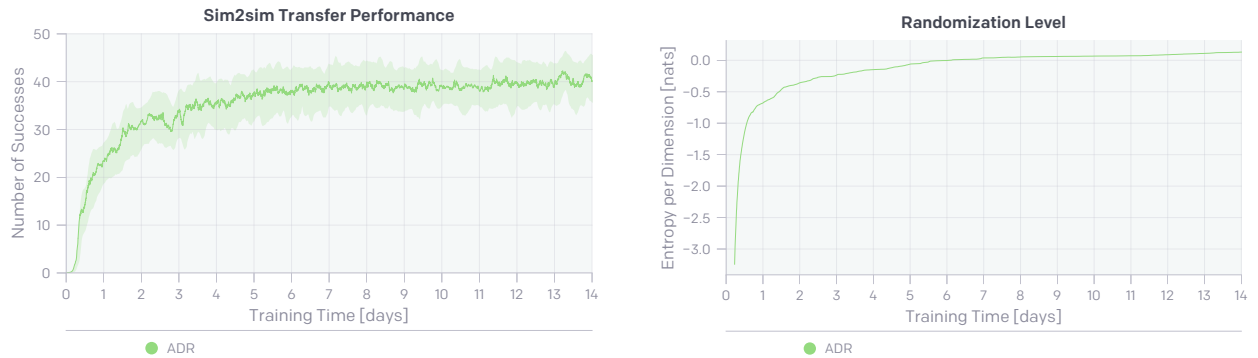


Figure 15: Sim2sim performance (left) and ADR entropy (right) over the course of training. We can see that a policy trained with ADR has better sim2sim transfer as ADR increases the randomization level over time.

We first consider the sim2sim case. More concretely, we train a policy with ADR and continuously benchmark its performance on an distribution of environments with manually tuned randomizations, very similar to the ones we used in [79]. Note that no ADR experiment has ever been trained on this distribution directly. Instead we use ADR to decide what distribution to train on, making the manually designed distribution over environments a test set for sim2sim transfer. We report our results in Figure 15.

As seen in Figure 15, the policy trained with ADR transfers to the manually randomized distribution. Furthermore, the sim2sim transfer performance increases as ADR increases the randomization entropy.

### 8.1.2 Sim2Real

Next, we evaluate the sim2real transfer capabilities of our policies. Since rollouts on the robot are expensive, we limit ourselves to 7 different policies that we evaluate. For each of them, we collect a total of 10 trials on the robot and measure the number of consecutive successes. As before, a trial ends when we either achieve 50 successes, the robot times out or the block is dropped. For each policy we deploy, we also report simulation performance by measuring the number of successes across 500 trials each for reference. As before, we use the manually designed randomizations as described in [79] for sim evaluations. We summarize our results in Table 3 and report detailed results in Appendix D.

Table 3: Performance of different policies on the block reorientation task. We evaluate each policy in simulation (N=500 trials) and on the real robot (N=10 trials) and report the mean  $\pm$  standard error and median number of successes. For ADR policies, we report the entropy in nats per dimension (npd). For “Manual DR”, we obtain an upper bound on its ADR entropy by running ADR with the policy fixed and report the entropy once the distribution stops changing (marked with an “\*”).

Policy	Training Time	ADR Entropy	Successes (Sim)		Successes (Real)	
			Mean	Median	Mean	Median
Baseline (data from [79])	—	—	$43.4 \pm 0.6$	<b>50</b>	$18.8 \pm 5.4$	13.0
Baseline (re-run of [79])	—	—	$33.8 \pm 0.9$	<b>50</b>	$4.0 \pm 1.7$	2.0
Manual DR	13.78 days	$-0.348^*$ npd	$42.5 \pm 0.7$	<b>50</b>	$2.7 \pm 1.1$	1.0
ADR (Small)	0.64 days	$-0.881$ npd	$21.0 \pm 0.8$	15	$1.4 \pm 0.9$	0.5
ADR (Medium)	4.37 days	$-0.135$ npd	$34.4 \pm 0.9$	<b>50</b>	$3.2 \pm 1.2$	2.0
ADR (Large)	13.76 days	$0.126$ npd	$40.5 \pm 0.7$	<b>50</b>	$13.3 \pm 3.6$	11.5
ADR (XL)	—	$0.305$ npd	$45.0 \pm 0.6$	<b>50</b>	$16.0 \pm 4.0$	12.5
ADR (XXL)	—	<b>0.393 npd</b>	<b><math>46.7 \pm 0.5</math></b>	<b>50</b>	<b><math>32.0 \pm 6.4</math></b>	<b>42.0</b>

The first two rows connect our results in this work to the previous results reported in [79]. For convenience, we repeat the numbers reported in [79] in the first row. We also re-deploy the exact same policy we used back then on our setup today. We find that the same policy performs much worse today, presumably because both our physical setup and simulation have changed since (as described in Section 3 and Section 4, respectively).

The next section of the table compares a policy trained with ADR and a policy trained with manual domain randomization (denoted as “Manual DR”). Note that “Manual DR” uses the same randomizations as the baseline from [79] but is trained on our current setup with the same model architecture and hyperparameters as the ADR policy. For the ADR policy, we select snapshots at different points during training at varying levels of entropy and denote them as small, medium, and large.<sup>13</sup> We can clearly see a pattern: increased ADR entropy corresponds to increased sim2sim and sim2real transfer. The policy trained with manual domain randomization achieves high performance in simulation. However, when deployed on the robot, it fails. This is because, in contrast to our results obtained in [79], we did not tune our simulation and randomization setup by hand to match changes in hardware. Our ADR policies transfer because ADR automates this process and results in training distributions that are vastly broader than our manually tuned distribution was in the past. Also note that “ADR (Large)” and “Manual DR” were trained for the same amount of wall-clock time and share all training hyperparameters except for the environment distribution, i.e. they are fully comparable. Due to compute constraints, we train those policies at  $1/4$ th of our usual scale in terms of compute (compare Section 6).

The last block of the table lists results that we obtained when scaling ADR up. We report results for “ADR (XL)” and “ADR (XXL)”, referring to two long-running experiments that were continuously trained for extended periods of time and at larger scale. We can see that they exhibit the best sim2sim and sim2real transfer and that, again, an increase in ADR entropy corresponds to vastly improved sim2real transfer. Our best result significantly beat the baseline reported in [79] even though we did not tune the simulation and robot setup for peak performance on the block reorientation task: we increase mean performance by almost  $2\times$  and median performance by more than  $3\times$ . As a side note, we also see that policies trained with ADR eventually achieve near-perfect performance for sim2sim transfer as well.

<sup>13</sup>Note that this is one experiment, not multiple different experiments, taken at different points in time during training.

In summary, ADR clearly leads to improved transfer with much less need for hand-engineered randomizations. We significantly outperformed our previous best results, which were the result of multiple months of iterative manual tuning.

## 8.2 Effect of Curriculum on Policy Training

We designed ADR to expand the complexity of the training distribution gradually. This makes intuitive sense: start with a single environment and then grow the distribution over environments as the agent progresses. The resulting curriculum should make it possible to eventually master a highly diverse set of environments. However, it is not clear if this curriculum property is important or if we can train with a fixed set of domain randomization parameters once they have been found.

To test for this, we conduct the following experiment. We train one policy with ADR on the block reorientation task and compare it against multiple policies with different fixed randomizations. We use 4 different fixed levels: small, medium, large, and XL. They correspond to the ADR parameters of the policies from the previous section (compare Table 3). However, note that we only use the ADR parameters, *not* the policies from Table 3. Instead, we train new policies from scratch using these parameters and train all of them for the same amount of wall-clock time. We evaluate performance of all policies continuously on the same manually randomized distribution from [79], i.e. we test for sim2sim transfer in all cases. We depict our results in Figure 16. Note that for all DR runs the randomization entropy is constant; only the one for ADR gradually increases.

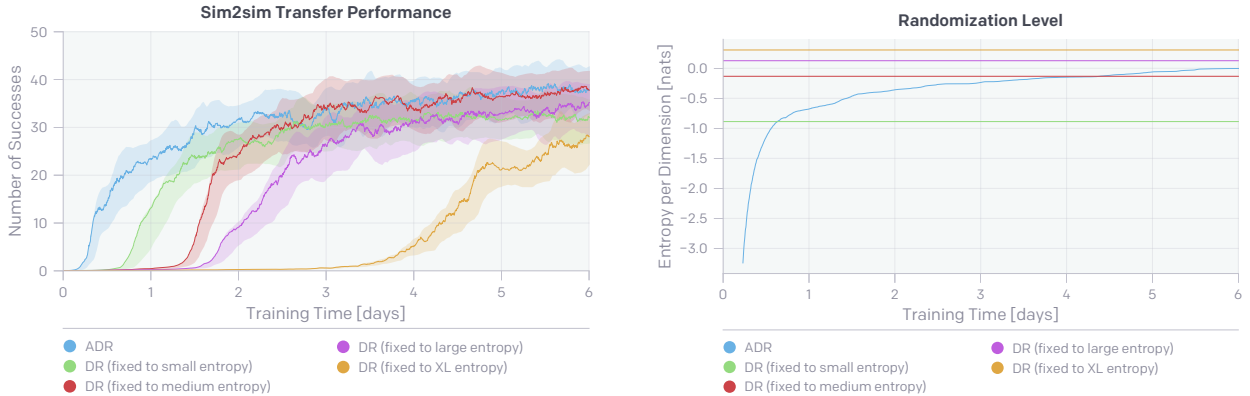


Figure 16: Sim2sim performance (left) and ADR entropy (right) over the course of training. *ADR* refers to a regular training run, i.e. we start with zero randomization and let ADR gradually expand the randomization level. We compare ADR against runs with domain randomization (DR) fixed at different levels and train policies on each of those environment distributions. We can see that ADR makes progress much faster due to its curriculum property.

Our results in Figure 16 clearly demonstrate that adaptively increasing the randomization entropy is important: the ADR run achieves high sim2sim transfer much more quickly than all other runs with fixed randomization entropy. There is also a clear pattern: the larger the fixed randomization entropy, the longer it takes to train from scratch. We hypothesize that for a sufficiently difficult task and randomization entropy, training from scratch becomes infeasible altogether. More concretely, we believe that for too complex environments the policy would never learn due to the task being so hard that there is no sufficient reinforcement learning signal.

## 8.3 Effect of ADR on Vision Model Performance

When training vision models, ADR controls both the ranges of randomization in ORRB (i.e. light distance, material metallic and glossiness) and TensorFlow distortion operations (i.e. adding Gaussian noise and channel noise). A full list of vision ADR randomization parameters are available in Section B.5 Table 11. We train ADR-enhanced vision models to do state estimation for both the block reorientation [79] and Rubik’s cube task. As shown in Table 4, we are able to reduce the prediction errors on both block orientation and position further than our manual domain randomization

<sup>13</sup>The attentive reader will notice that the sim2sim performance reported in Figure 16 is different from the sim2sim performance reported in Table 3. This is because here we train the policies for *longer* but on a *fixed* ADR entropy whereas in Table 3 we had a single ADR run and took snapshots at different points over the course of training.

results in the previous work [79].<sup>14</sup> We can also see that increased ADR entropy again corresponds to better sim2real transfer.

Table 4: Performance of vision models at different ADR entropy levels for the block reorientation state estimation task. Note that the baseline model here uses the same manual domain randomization configuration as in [79] but is evaluated on a newly collected real image dataset (the same real dataset described in Table 2)

Model	Training Time	ADR Entropy	Errors (Sim)		Errors (Real)	
			Orientation	Position	Orientation	Position
Manual DR	13.62 hrs	—	1.99°	4.03 mm	5.19°	8.53 mm
ADR (Small)	2.5 hrs	0.922 npd	2.81°	4.21 mm	6.99°	8.13 mm
ADR (Middle)	3.87 hrs	1.151 npd	2.73°	4.11 mm	6.66°	8.14 mm
ADR (Large)	12.76 hrs	<b>1.420 npd</b>	1.85°	5.18 mm	<b>5.09°</b>	<b>7.85 mm</b>

Predicting the full state of a Rubik’s cube is a more difficult task and demands longer training time. A vision model using ADR succeeds to achieve lower errors than the baseline model with manual DR configuration given a similar amount of training time, as demonstrated in Table 5. Higher ADR entropy well correlates with lower errors on real images. ADR again outperforms our manually tuned randomizations (i.e. the baseline). Note that errors in simulation increase as ADR generates harder and harder synthetic tasks.

Table 5: Performance of vision models at different ADR entropy levels for the Rubik’s cube prediction task (See Section 7). The real image datasets used for evaluation are same as in Table 2. The full evaluation results are reported in Appendix D Table 22.

Model	Training Time	ADR Entropy	Errors (Sim)			Errors (Real)		
			Orientation	Position	Top angle	Orientation	Position	Top angle
Baseline	76.29 hrs	—	6.52°	2.63 mm	11.95°	7.81°	6.47 mm	15.92°
ADR (Small)	20.9 hrs	−0.565 npd	5.02°	3.36 mm	9.34°	8.93°	7.61 mm	16.57°
ADR (Middle)	30.6 hrs	0.511 npd	15.68°	3.02 mm	20.29°	8.44°	7.30 mm	15.81°
ADR (Large)	75.1 hrs	<b>0.806 npd</b>	15.76°	3.58 mm	20.78°	<b>7.48°</b>	<b>6.24 mm</b>	<b>13.83°</b>

## 8.4 Solving the Rubik’s Cube

In this section, we push the limits of sim2real transfer by considering a manipulation problem of unprecedented complexity: solving Rubik’s cube using the real Shadow hand. This is a daunting task due to the complexity of Rubik’s cube and the interactions between it and the hand: in contrast to the block reorientation task, there is no way we can accurately capture the object in simulation. While we model the Rubik’s cube (see Section 4), we make no effort to calibrate its dynamics. Instead, we use ADR to automate the randomization of environments.

We further need to sense the state of the Rubik’s cube, which is also much more complicated than for the block reorientation task. We always use vision for the pose estimation of the cube itself. For the 6 face angles, we experiment with two different setups: the Giiker cube (see Section 3) and a vision model which predicts face angles (see Section 7).

We first evaluate performance on this task quantitatively and then highlight some qualitative findings.

### 8.4.1 Quantitative Results

We compare four different policies: a policy trained with manual domain randomization (“Manual DR”) using the randomizations that we used in [79] trained for about 2 weeks, a policy trained with ADR for about 2 weeks, and two policies we continuously trained and updated with ADR over the course of months.

<sup>14</sup>Note that this model has the same manual DR as in [79] but is evaluated on a newly collected real image set, so the numbers are slightly different from [79].

Table 6: Performance of different policies on the Rubik’s cube for a fixed fair scramble goal sequence. We evaluate each policy on the real robot (N=10 trials) and report the mean  $\pm$  standard error and median number of successes (meaning the total number of successful rotations and flips). We also report two success rates for applying half of a fair scramble (“half”) and the other one for fully applying it (“full”). For ADR policies, we report the entropy in nats per dimension (npd). For “Manual DR”, we obtain an upper bound on its ADR entropy by running ADR with the policy fixed and report the entropy once the distribution stops changing (marked with an “\*”).

Policy	Sensing		ADR Entropy	Successes (Real)		Success Rate	
	Pose	Face Angles		Mean	Median	Half	Full
Manual DR	Vision	Giiker	−0.569* npd	$1.8 \pm 0.4$	2.0	0 %	0 %
ADR	Vision	Giiker	−0.084 npd	$3.8 \pm 1.0$	3.0	0 %	0 %
ADR (XL)	Vision	Giiker	0.467 npd	$17.8 \pm 4.2$	12.5	30 %	10 %
ADR (XXL)	Vision	Giiker	<b>0.479 npd</b>	<b><math>26.8 \pm 4.9</math></b>	<b>22.0</b>	<b>60 %</b>	<b>20 %</b>
ADR (XXL)	Vision	Vision	<b>0.479 npd</b>	$12.8 \pm 3.4$	10.5	20 %	0 %

To evaluate performance, we define a fixed procedure that we repeat 10 times per policy to obtain 10 trials. More concretely, we always start from a solved cube state and ask the hand to move the Rubik’s cube into a fair scramble. Since the problem is symmetric, this is equivalent to solving the Rubik’s cube starting from a fairly scrambled Rubik’s cube. However, it reduces the probability of human error and labor significantly since ensuring a correct initial state is much simpler if the cube is solved. We use the following fixed scrambling sequence for all 10 trials, which we obtained using the “TNoodle” application of the World Cube Association<sup>15</sup> via a random sample (i.e., this was not cherry-picked):

L2 U2 R2 B D2 B2 D2 L2 F’ D’ R B F L U’ F D’ L2

Completing this sequence requires a total of 43 successes (26 face rotations and 17 cube flips). If the sequence is completed successfully, we continue the trial by reversing the sequence. A trial ends if 50 successes have been achieved, if the cube is dropped, or if the policy fails to achieve a goal within 1600 timesteps, which corresponds to 128 seconds.

For each trial, we measure the number of successfully achieved goals (both flips and rotations). We also define two thresholds for each trial: Applying at least half of the fair scramble successfully (i.e. 22 successes) and applying at least the full fair scramble successfully (i.e. 43 successes). We report success rates for both averaged across all 10 trials and denote them as “half” and “full”, respectively. Achieving the “full” threshold is equivalent to solving the Rubik’s cube since going from solved to scrambled is as difficult as going from scrambled to solved.<sup>16</sup> We report our summarized results in Table 6 and full results in Appendix D.

We see a very similar pattern as before: manual domain randomization fails to transfer. For policies that we trained with ADR, we see that sim2real transfer clearly depends on the entropy per dimension. “Manual DR” and “ADR” were trained for 14 days at  $1/4$ th of our usual scale in terms of compute (see Section 6) and are fully comparable. Our best policy, which was continuously trained over the course of multiple months at larger scale, achieves 26.80 successes on average over 10 trials. This corresponds to successfully solving a Rubik’s cube that requires 15 face rotations 60% of the time and to solve a Rubik’s cube that requires 26 face rotations 20% of the time. Note that 26 quarter face rotations is the worst case for solving a Rubik’s cube with only about 3 Rubik’s cube configurations requiring that many [88]. In other words, almost all solution sequences will require less than 26 face rotations.

### 8.4.2 Qualitative Results

We observe many interesting emergent behaviors on our robot when using our best policy (“ADR (XXL)”) for solving the Rubik’s cube. We encourage the reader to watch the uncut video footage we recorded: <https://youtu.be/kVmp0uGtShk>.

<sup>15</sup><https://www.worldcubeassociation.org/regulations/scrambles/>

<sup>16</sup>With the exception of the fully solved configuration being slightly easier for the vision model to track.

<sup>16</sup>This video solves the Rubik’s cube from an initial randomly scrambled state. This is different from the quantitative experiments we conducted in the previous section.



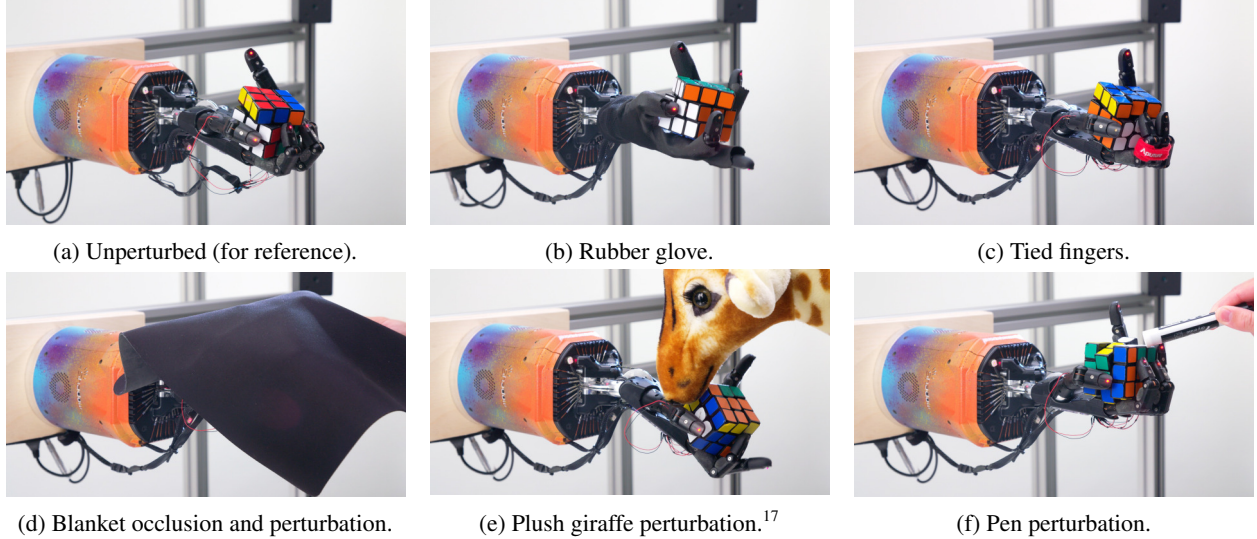


Figure 17: Example perturbations that we apply to the real robot hand while it solves the Rubik’s cube. We did not train the policy to be able to handle those perturbations, yet we observe that it is robust to all of them. A video of is available: <https://youtu.be/QyJGxc9WeNo>

For example, we observe that the robot sometimes accidentally rotates an incorrect face. If it does so, our best policies are usually able to recover from this mistake by first rotating the face back and then pursuing the original subgoal without us having to change the subgoal. We also observe that the robot first aligns faces after performing a flip before attempting a rotation to avoid interlocking due to misalignment. Still, rotating a face can be challenging at times and we sometimes observe situations in which the robot is stuck. In this case we often see that the policy eventually adjusts its grasp to attempt the face rotation a different way, thus often succeeding eventually. Other times we observe our policy attempting a face rotation but the cube slips, resulting in a rotation of the entire cube as opposed to a specific face. In this case the policy rearranges its grasp and tries again, usually succeeding eventually.

We also observe that the policy appears more likely to drop the cube after being stuck on a challenging face rotation for a while. We do not quantify this but hypothesize that it might have “forgotten” about flips by then since the recurrent state of the policy has only observed a mostly stationary cube for several seconds. For flips, information about the cube’s dynamics properties are more important. Similarly, we also observe that the policy appears to be more likely to drop the cube early on, presumably again because the necessary information about the cube’s dynamic properties have not yet been captured in the policy’s hidden state.

We also experiment with several perturbations. For example, we use a rubber glove to significantly change the friction and surface geometry of the hand. We use straps to tie together multiple fingers. We use a blanket to occlude the hand and Rubik’s cube during execution. We use a pen and plush giraffe to poke the cube. While we do not quantify these experiments, we find that our policy still is able to perform multiple face rotations and cube flips under all of these conditions even though it was clearly not trained on them. Figure 17 shows examples of perturbations we tried. A video showing the behavior of our policy under those perturbations is also available: <https://youtu.be/QyJGxc9WeNo>

## 9 Signs of Meta-Learning

We believe that a sufficiently diverse set of environments combined with a memory-augmented policy like an LSTM leads to *emergent meta-learning*. In this section, we systematically study our policies trained with ADR for signs of meta-learning.

### 9.1 Definition of Meta-Learning

Since we train each policy on only one specific task (i.e. the block reorientation task or solving the Rubik’s cube), we define meta-learning in our context as learning about the dynamics of the underlying Markov decision process.

<sup>17</sup>Her name is Rubik, for obvious reasons.

More concretely, we are looking for signs where our policy updates its belief about the true transition probability  $P(s_{t+1} | s_t, a_t)$  as it observes data over time.

In other words, when we say “meta-learning”, what we really mean is learning to learn about the environment dynamics. Within other communities, this is also called on-line system identification. In our case though, this is an emergent property.

## 9.2 Response to Perturbations

We start by studying the behavior of our policy and how it responds to a variety of perturbations to the dynamics of the environment. We conduct all experiments in simulation and use the Rubik’s cube task. In all our experiments, we fix the type of subgoal we consider to be either cube flips or cube rotations. We run the policy until it achieves the 10<sup>th</sup> flip (or rotation) and then apply a perturbation. We then continue until the 30<sup>th</sup> successful flip (or rotation) and apply another perturbation. We measure the time it took to achieve the 1<sup>st</sup>, 2<sup>nd</sup>, . . . , 50<sup>th</sup> flip (or rotation) for each trial, which we call “time to completion”. We also measure during which flip (or rotation) the policy failed. By averaging over multiple trials that we all run in simulation, we can compute the average time to completion and failure probability *per flip* (or rotation).

If our policy learns at test time, we’d expect the average time to completion to gradually decrease as the policy learns to identify the dynamics of its concrete environment and becomes more efficient as it accumulates more information over the course of a trial. Once we perturb the system, however, the policy needs to update its belief. We therefore expect to see a spike, i.e. achieving a flip (or rotation) should take longer after a perturbation but should again decrease as the policy readjusts its belief about the perturbed environment. Similarly, we expect to see the failure probability to be higher during the first few flips (or rotations) since the policy has had less time to learn about its environment. We also expect the failure probability to spike after each perturbation.

We experiment with the following perturbations:

- **Resetting the hidden state.** During a trial, we reset the hidden state of the policy. This leaves the environment dynamics unchanged but requires the policy to re-learn them since its memory has been wiped.
- **Re-sampling environment dynamics.** This corresponds to an abrupt change of environment dynamics by resampling the parameters of all randomizations while leaving the simulation state<sup>18</sup> and hidden state intact.
- **Breaking a random joint.** This corresponds to us disabling a randomly sampled joint of the robot hand by preventing it from moving. This is a more nuanced experiment since the overall environment dynamics are the same but the way in which the robot can interact with the environment has changed.

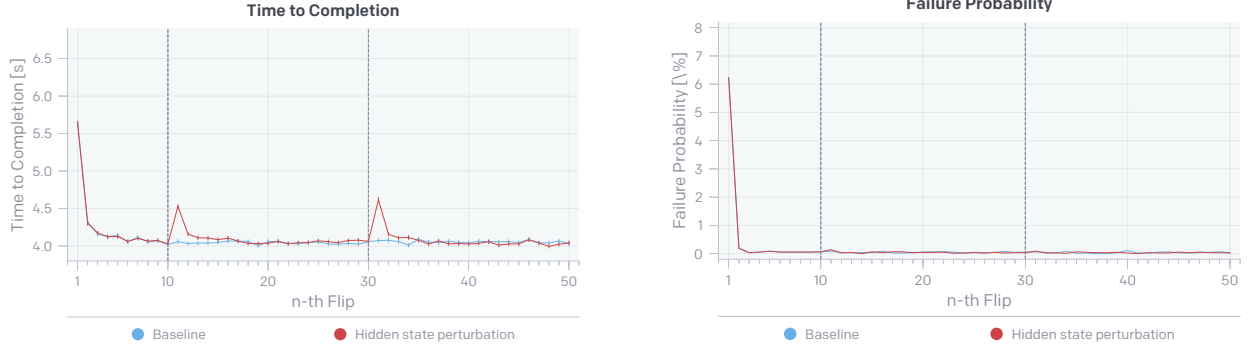
We show the results of our experiments for cube flips in Figure 18. The same plot is available for cube face rotations in Section D.4. For the average time to completion, we only include trials that achieved 50 flips to avoid inflating our results.<sup>19</sup>

Our results show a very clear trend: for all runs, we observe a clear adjustment period over the first few cube flips. Achieving the first one takes the longest with subsequent flips being achieved more and more quickly. Eventually the policy converges to approximately 4 seconds per flip on average, which is an improvement of roughly 1.6 seconds compared to the first flip. This was exactly what we predicted: if the policy truly learns at test time by updating its recurrent state, we would expect it to become gradually more efficient. The same holds true for failure probabilities: the policy is much more likely to fail early.

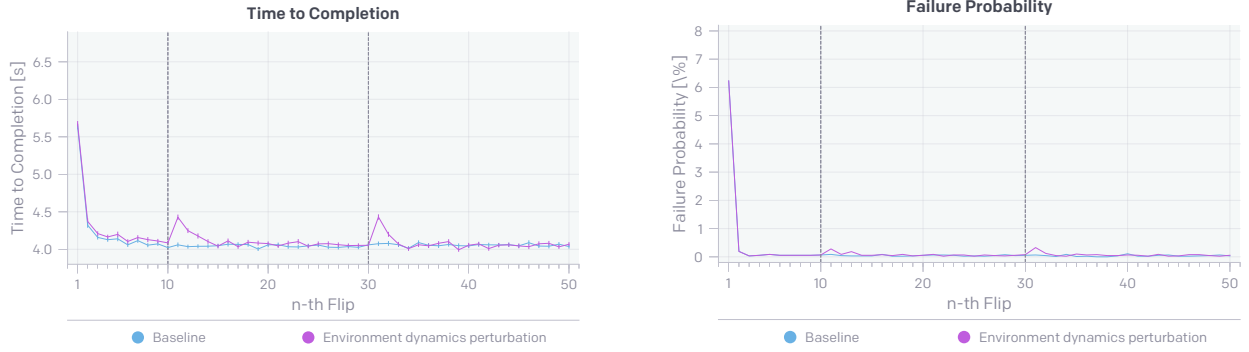
When we reset the hidden state of our policy (compare Figure 18a), we can see the time to completion spike up significantly immediately after. This is because the policy again needs to identify the environment since all its memory has been wiped. Note that the spike in time to completion is much less than the initial time to completion. This is the case because we randomize the initial cube position and configuration. In contrast, when we reset the hidden state, the hand had manipulated the cube before so it is in a more beneficial position for the hand to continue after the hidden state is reset. This is also visible in the failure probability, which is close to zero even after the perturbation is applied, again because the cube is already in a beneficial position which makes it less likely to be dropped.

<sup>18</sup>The simulation state is the current kinematic configuration of the cube, the hand, and the goal.

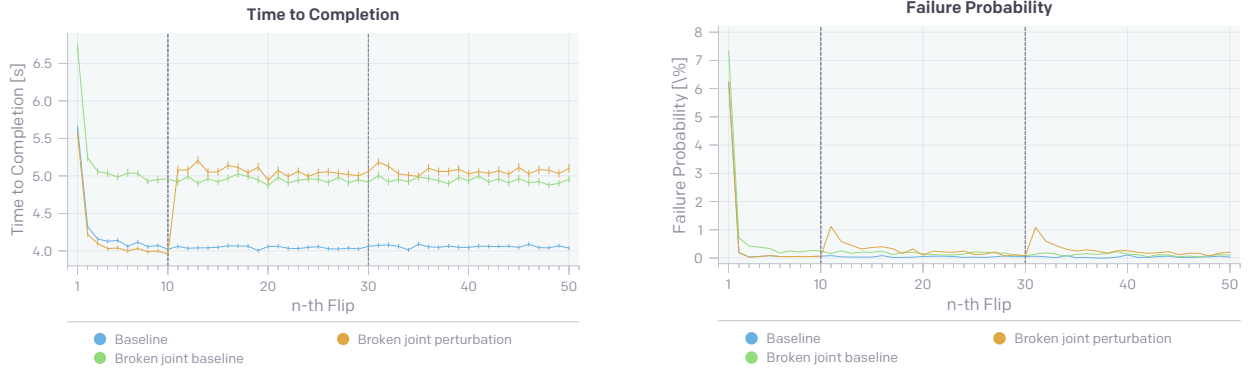
<sup>19</sup>To explain further: By only include trials with 50 successes, we ensure that we measure performance over a *fixed* distribution over environments, and thus only measure how the performance of the policy changes across this fixed set. If we would not restrict this, harder environments would be more likely to lead to failures within the first few successes after a perturbation and then would “drop out”. The remaining ones would be easier and thus even a policy without the ability to adapt would appear to improve in performance. Failures are still important, of course, which is why we report them in the failure probability plots.



(a) Resetting the hidden state.



(b) Re-sampling environment dynamics.



(c) Breaking a random joint.

Figure 18: We run 10 000 simulated trials with only cube flips until 50 flips have been achieved. For each of the cube flips (i.e. the 1<sup>st</sup>, 2<sup>nd</sup>, ..., 50<sup>th</sup>), we measure the average time to completion (in seconds) and average failure probability over those 10k trials. Error bars indicate the estimated standard error. “Baseline” refers to a run without any perturbations applied. “Broken joint baseline” refers to trials where a joint was randomly disabled from the very beginning. We then compare against trials that start without any perturbations but are perturbed at the marked points after the 10<sup>th</sup> and 30<sup>th</sup> flip by (a) resetting the policy hidden state, (b) re-sampling environment dynamics, or (c) breaking a random joint.

The second experiment perturbs the environment by resetting its environment dynamics but keeping the simulation state itself intact (see Figure 18b). We see a similar effect as before: after the perturbation is applied, the time to completion spikes up and then decreases again as the policy adjusts. Interestingly this time the policy is more likely to fail compared to resetting the hidden state. This is likely the case because the policy is “surprised” by the sudden change and performed actions that would have been appropriate for the old environment dynamics but led to failure in the new ones.

An especially interesting experiment is the broken joint one (see Figure 18c): for the “broken joint baseline”, we can see how the policy adjusts over long time horizons, with improvements clearly visible over the course of all 50 flips in both time to completion and failure probability. In contrast, “broken joint perturbation” starts with all joints intact. After the perturbation, which breaks a random joint, we see a significant jump in the failure probability, which then gradually decreases again as the policy learns about this limitation. We also find that the “broken joint perturbation” performance never quite catches up to the “Broken joint baseline”. We hypothesize that this could be because the policy has already “locked-in” some information in its recurrent state and therefore is not as adjustable anymore. Alternatively, maybe it just has not accumulated enough information yet and the baseline policy is in the lead because it has an “information advantage” of at least 10 achieved flips. We found it very interesting that our policy can learn to adapt internally to broken joints. This is in contrast to prior work that explicitly searched over a set of policies until it found one that works for a broken robot [26]. Note however that the “Broken joint baseline” never fully matches the performance of the “Baseline”, suggesting that the policy can not fully recover performance.

In summary, we find clear evidence of our policy learning about environment dynamics and adjusting its behavior accordingly to become more efficient at *test time*. All of this learning is emergent and only happens by updating the policy’s recurrent state.

### 9.3 Recurrent State Analysis

We conducted experiments to study whether the policy has learned to infer and store useful information about the environment in its recurrent state. We consider such information to be strong evidence of meta-learning, since no explicit information regarding the environment parameters was provided during training time.

The main method that we use to probe the amount of useful environment information is to predict environment parameters, such as cube size or the gravitational constant, from the policy LSTM hidden and cell states. Given a policy with an LSTM with hidden states  $\mathbf{h}$  and cell states  $\mathbf{c}$ , we use  $\mathbf{z} = \mathbf{h} + \mathbf{c}$  as the prediction model input. For environment parameter  $p$ , we trained a simple prediction model  $f_p(\mathbf{z})$  containing one hidden layer with 64 units and ReLU activation, followed by a sigmoid output. The outputs correspond to the probability that the value of  $p$  for the environment is greater or smaller than the average randomized value.

The prediction model was trained on hidden states collected from policy rollouts at time step  $t$  in a set of environments  $\mathcal{E}_t$ . Each environment  $e \in \mathcal{E}_t$  contains a different value of the parameter  $p$ , sampled according to its randomization range. To observe the change in the stored information over time, we collected data from times steps  $t \in \{1, 30, 60, 120\}$  where each time step is equal to  $\Delta t = 0.08$  s in simulation. We used the cross-entropy loss and trained the model until convergence. The model was tested on a new set of environments  $\mathcal{F}_t$  with newly sampled values of  $p$ .

#### 9.3.1 Prediction Accuracy over Time

We studied four different environment parameters, listed in Table 7 along with their randomization ranges. The randomization ranges were used to generate the training and test environment sets  $\mathcal{E}_{t,p}$ ,  $\mathcal{F}_{t,p}$ . Each parameter is sampled using a uniform distribution over the given range. Randomization ranges were taken from the end of ADR training.

Table 7: Prediction environment parameters and their randomization ranges in physical units. We predict whether or not the parameter is larger or smaller than the given average.

Parameter	Block Reorientation		Rubik’s Cube	
	Average	Range	Average	Range
Cube size [m]	0.055	[0.046, 0.064]	0.057	[0.049, 0.066]
Time step [s]	0.1	[0.05, 0.15]	0.09	[0.05, 0.13]
Gravity [ $\text{m s}^{-2}$ ]	9.80	[6.00, 14.0]	9.80	[6.00, 14.0]
Cube mass [kg]	0.0780	[0.0230, 0.179]	0.0902	[0.0360, 0.158]

Figure 19 shows the test accuracy of the trained prediction models for block reorientation and Rubik’s cube policies. Observe that prediction accuracy at the start of a rollout is near random guessing, since no useful information is stored in the hidden states.

As rollouts progress further and the policy interacts with the environment, prediction accuracy rapidly improves to over 80% for certain parameters. This is evidence that the policy is successfully inferring and storing useful information

regarding the environment parameters in its LSTM hidden and cell states. Note that we do not train the policy explicitly to store information about those semantically meaningful physical parameters.

There exists some variability in the prediction accuracy over different parameters and between block reorientation and Rubik’s cube policies. For example, note that the prediction accuracy for cube size (over 80%) is consistently higher than that of cube mass (50 – 60%). This may be due to the relative importance of cube size and mass to policy performance; i.e., a heavier cube changes the difficulty of Rubik’s cube face rotation less than a larger cube. There also exist differences for the same parameter between tasks: For the cube mass, the block reorientation policy stores more information about it than the Rubik’s cube policy. We hypothesize that this is because the block reorientation policy uses a dynamic approach that tosses the block around to flip it. In contrast, the Rubik’s cube policy flips the cube much more deliberately in order to avoid unintentional misalignments of the cube faces. For the former, knowing the cube mass is therefore more important since the policy needs to be careful to not apply too much force. We believe the variations in prediction accuracy for the other parameters and the two policies also reflect the relative importance of each parameter to the policy and the given task.

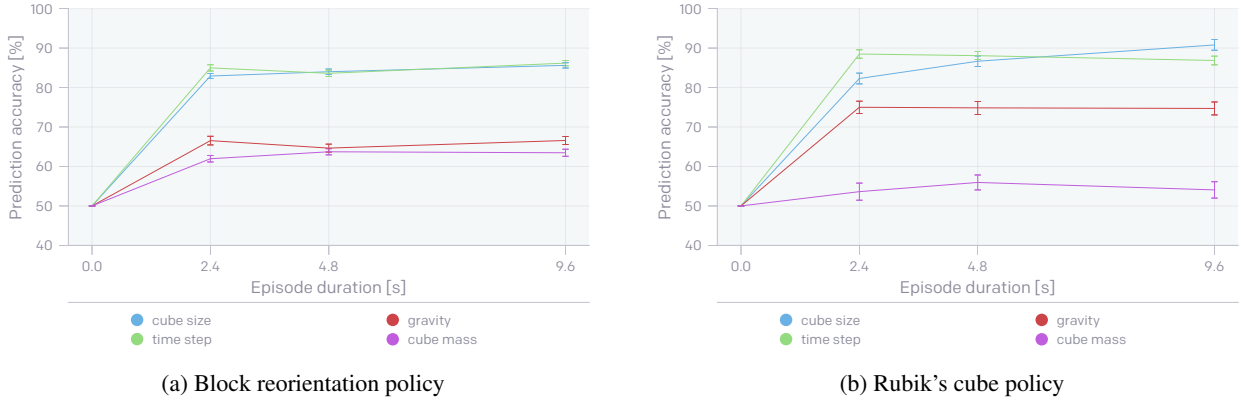


Figure 19: Test accuracy of environment parameter prediction model based on the hidden states of (a) block reorientation and (b) Rubik’s cube policies. Error bars denote the standard error.

### 9.3.2 Information Gain over Time

To further study the information contained in a policy hidden state and how it evolves during a rollout, we expanded the prediction model in the above experiments to predict a set of 8 equally-spaced discretized values (“bins”) within a parameter’s randomization range. We consider the output probability distribution of the predictor to be an approximate representation of the posterior distribution over the environment parameter, as inferred by the policy.

In Figure 20, we plot the entropy of the predictor output distribution in test environments over rollout time. The parameter being predicted is cube size. The results clearly show that the posterior distribution over the environment parameters converges to a certain final distribution as a rollout progresses. The convergence speed is rather fast at below 5.0 seconds. Notice that this is consistent with our perturbation experiments (compare Figure 19): the first flip roughly takes 5 – 6 seconds and we see a significant speed-up after. Within this time, the information gain for the cube size parameter is approximately 0.9 bits. Interestingly, the entropy eventually seems to converge to 2.0 bits and then stops decreasing. This again highlights that our policies only store the amount of information they need to act optimally.

### 9.3.3 Prediction Accuracy and ADR Entropy

We performed the hidden state prediction experiments for block reorientation policies trained using ADR with different values of ADR entropy. Since we believe that the ability of the policy to infer and store (i.e., meta-learn) useful information regarding environment parameters is correlated with the diversity of the environments used during training, we expect the prediction accuracy to be positively correlated with the policy’s ADR entropy.

Four block reorientation policies corresponding to increasing ADR entropy were used in the following experiments. We seek to predict the cube size parameter at 60 rollout time steps, which corresponds to 4.8 seconds of simulated time. The test accuracies are shown in Table 8. The results indicate that prediction accuracy (hence information stored in hidden states) is strongly correlated with ADR entropy.



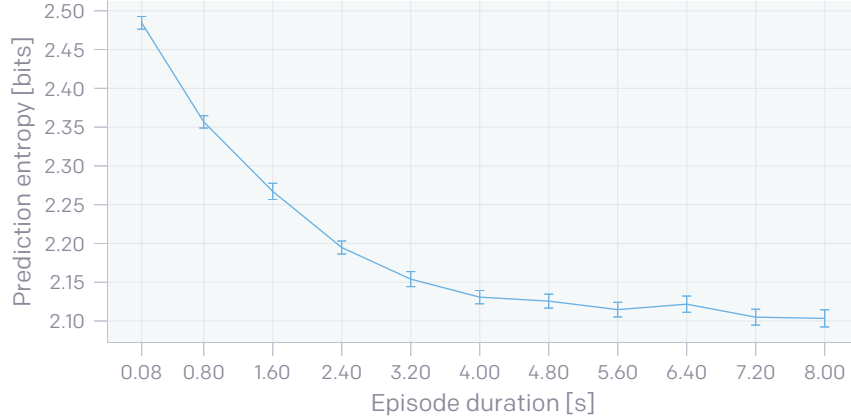


Figure 20: Mean prediction entropy over rollout time for an 8-bin output predictor. Error bars denote the standard error. The predictor was trained at a fixed 4.8 seconds during rollouts. Parameter being predicted is cube size. Note the information gain of 0.9 bits in less than 5.0 seconds. For reference, the entropy for random guessing (i.e. uniform probability mass over all 8 bins) is 3 bits.

Table 8: Block reorientation policy hidden state prediction over ADR entropy. The environment parameter predicted is cube size. ADR entropy is defined to be nats per environment parameterization dimension or npd. All predictions were performed at rollout time step  $t = 60$  (which corresponds to 4.8 seconds of simulated time). We report mean prediction accuracy  $\pm$  standard error.

Policy	Training Time	ADR Entropy	Prediction Accuracy
ADR (Small)	0.64 days	$-0.881$ npd	$0.68 \pm 0.021$
ADR (Medium)	4.37 days	$-0.135$ npd	$0.75 \pm 0.027$
ADR (Large)	13.76 days	$0.126$ npd	$0.79 \pm 0.022$
ADR (XL)	—	$0.305$ npd	<b><math>0.83 \pm 0.014</math></b>

### 9.3.4 Recurrent state visualization

We used neural network interpretability techniques [17, 77] to visualize policy recurrent states during a rollout. We found distinctive activation patterns that correspond to high-level skills exhibited by the robotic hand. See Section E.1 for details.

## 10 Related Work

### 10.1 Dexterous Manipulation

Dexterous manipulation has been an active area of research for decades [35, 92, 13, 76, 66]. Many different approaches and strategies have been proposed over the years. This includes rolling [14, 43, 44, 19, 30], sliding [19, 101], finger gaiting [44], finger tracking [91], pushing [28], and re-grasping [110, 29]. For some hand types, strategies like pivoting [7], tilting [33], tumbling [97], tapping [48], two-point manipulation [6], and two-palm manipulation [32] are also options. These approaches use planning and therefore require exact models of both the hand and object. After computing a trajectory, the plan is typically executed open-loop, thus making these methods prone to failure if the model is not accurate.<sup>20</sup>

Other approaches take a closed-loop approach to dexterous manipulation and incorporate sensor feedback during execution, e.g. tactile sensing [105, 61, 62, 63]. While those approaches allow for the correction of mistakes at runtime, they still require reasonable models of the robot kinematics and dynamics, which can be challenging to obtain for under-actuated hands with many degrees of freedom.

<sup>20</sup>Some methods use iterative re-planning to partially mitigate this issue.

Deep reinforcement learning has also been used successfully to learn complex manipulation skills on physical robots. Guided policy search [58, 60] learns simple local policies directly on the robot and distills them into a global policy represented by a neural network. Soft Actor-Critic[42] has been recently proposed as a state-of-the-art model-free algorithm optimizing concurrently both expected reward and action entropy, that is capable of learning complex behaviors directly in the real world.

Alternative approaches include using many physical robots simultaneously, in order to be able to collect sufficient experience [38, 59, 51] or leveraging a model-based learning algorithms, which generally possess much more favorable sample complexity characteristics [121]. Some researchers have successfully utilized expert human demonstrations in guiding the training process of the agents [21, 65, 9, 122].

## 10.2 Dexterous In-Hand Manipulation

Since a very large body of past work on dexterous manipulation exists, we limit the more detailed discussion to setups that are most closely related to our work on dexterous in-hand manipulation.

Mordatch et al. [72] and Bai et al. [11] propose methods to generate trajectories for complex and dynamic in-hand manipulation, but their results are limited to simulation. There has also been significant progress in learning complex in-hand dexterous manipulation [85, 12], tool use [87] and even solving a smaller model of a Rubik’s Cube [64] using deep reinforcement learning, but those approaches were likewise only evaluated in simulation.

In contrast, multiple authors learn policies for dexterous in-hand manipulation directly on the robot. Hoof et al. [114] learn in-hand manipulation for a simple 3-fingered gripper whereas Kumar et al. [55, 54] and Falco et al. [34] learn such policies for more complex humanoid hands. In [73], the authors learn a forward dynamics model and use model predictive control to manipulate two Baoding balls with a Shadow hand. While learning directly on the robot means that modeling the system is not an issue, it also means that learning has to be performed with limited data. This is only possible when learning simple (e.g. linear or local) policies that, in turn, do not exhibit sophisticated behaviors.

## 10.3 Sim to Real Transfer

*Domain adaption* methods [112, 40], progressive nets [93], and learning inverse dynamics models [22] were all proposed to help with sim to real transfer. All of these methods assume access to real data. An alternative approach is to make the policy itself more adaptive during training in simulation using *domain randomization*. Domain randomization was used to transfer object pose estimators [107] and vision policies for fly drones [94]. This idea has also been extended to dynamics randomization [81, 8, 106, 119, 79] to learn a robust policy that transfers to similar environments but with different dynamics. Domain randomization was also used to plan robust grasps [67, 68, 108] and to transfer learned locomotion [106] and grasping [123] policies for relatively simple robots. Pinto et al. [84] propose to use *adversarial training* to obtain more robust policies and show that it also helps with transfer to physical robots [83]. Hwangbo et al. [50] used real data to learn an actuation model and combined it with domain randomization to successfully transfer locomotion policies.

A number of recent works have focused on adapting the environment distribution of domain randomization to improve sim-to-real transfer performance. For policy training, one approach viewed the problem as a bi-level optimization [115, 90]. Chebotar et al. [18] used real-world trajectories and a discrepancy metric to guide the distribution search. In [70], a discriminator was used to guide the distribution in simulation. For vision models, domain randomization has been modified to improve image content diversity [52, 86, 25] and to adapt the distribution by using an adversarial network [120].

## 10.4 Meta-Learning via Reinforcement Learning

Despite being a very young field, meta-learning in the context of deep reinforcement learning already has a large body of work published. Algorithms such as MAML [36] and SNAIL [71] have been developed to improve the sample efficiency of reinforcement learning agents. A common theme in research is to try to exploit a shared structure in a distribution of environments, to quickly identify and adapt to previously unseen cases [57, 95, 49, 56] There are works that directly treat meta-learning as identification of a dynamics model of the environment [23] while others tackle the problem of task discovery for training [41]. Meta-learning has also been studied in a multi-agent setting [80].

The approach we’ve taken is directly based on  $RL^2$  [31, 116] where a general-purpose optimization algorithm trains a model augmented with memory to perform independent learning algorithm in the inner loop. The novelty in our results comes from the combination of automated curriculum generation (ADR), a challenging underlying problem (solving Rubik’s Cube) and a completely out-of-distribution test environment (sim2real).

This approach coincides with what Jeff Clune described as AI-GA [24], the AI-generating algorithms, whose one of three pillars is generating effective and diverse learning environments. In similar spirit to our ADR, related works include PowerPlay [98, 104], POET [117] and different varieties of self-play [102, 103, 78].

### 10.5 Robotic Systems Solving Rubik’s Cube

While many robots capable of solving a Rubik’s Cube exist today [2, 4, 3, 46], all of them which we are aware of have been built exclusively for this purpose and cannot generalize to other manipulation tasks.

## 11 Conclusion

In this work, we introduce automatic domain randomization (ADR), a powerful algorithm for sim2real transfer. We show that ADR leads to improvements over previously established baselines that use manual domain randomization for both vision and control. We further demonstrate that ADR, when combined with our custom robot platform, allows us to successfully solve a manipulation problem of unprecedented complexity: solve a Rubik’s cube using a real humanoid robot, the Shadow Dexterous Hand. By systematically studying the behavior of our learned policies, we find clear signs of emergent meta-learning. Policies trained with ADR are able to adapt at deployment time to the physical reality, which it has never seen during training, via updates to their recurrent state.

## Acknowledgements

We would like to thank Shadow Robot Company Ltd. for building, maintaining, and improving the Shadow Dexterous Hand, PhaseSpace Inc. for building custom Rubik’s cubes and supporting our motion capture setup, JITX for building custom electronics and writing firmware for the Giiker cube, and Kimly Do for designing the cage. We would also like to thank Chris Hallacy, Jacob Hilton, and Ludwig Schubert for volunteering as human Rubik’s cube solvers and everybody at OpenAI for their help and support.

We would also like to thank the following people for providing thoughtful feedback on earlier versions of this manuscript: Josh Achiam, Nick Cammarata, Jeff Clune, Harri Edwards, David Farhi, Ken Goldberg, Lerrel Pinto, and John Schulman.

## Author Contributions

This manuscript is the result of the work of the entire OpenAI Robotics team. We list the contributions of every team member here grouped by topic and in alphabetical order.

- Marcin Andrychowicz, Matthias Plappert, and Raphael Ribas developed the first version of ADR.
- Maciek Chociej, Alex Paino, Peter Welinder, Lilian Weng, and Qiming Yuan developed the vision state estimator.
- Maciek Chociej, Peter Welinder and Lilian Weng applied ADR to vision.
- Ilge Akkaya, Marcin Andrychowicz, Matthias Plappert, Raphael Ribas, Nikolas Tezak, Jerry Tworek, and Lei Zhang contributed to the RL training setup.
- Mateusz Litwin, Jerry Tworek, and Lei Zhang improved ADR for RL training.
- Mateusz Litwin improved the robot model and Jerry Tworek developed the Rubik’s cube simulation model.
- Arthur Petron and Jonas Schneider built the physical robot setup.
- Ilge Akkaya, Maciek Chociej, Mateusz Litwin, Arthur Petron, Glenn Powell, Jonas Schneider, Peter Welinder, Lilian Weng, and Qiming Yuan contributed software for the robot platform.
- Mateusz Litwin, Glenn Powell, and Jonas Schneider developed system infrastructure.
- Ilge Akkaya, Mateusz Litwin, Arthur Petron, Alex Paino, Matthias Plappert, Jerry Tworek, Lilian Weng, and Lei Zhang wrote this manuscript.
- Marcin Andrychowicz, Bob McGrew, Matthias Plappert, Jonas Schneider, Peter Welinder, and Wojciech Zaremba led aspects of this project and set research directions.
- Wojciech Zaremba led the team.

## References

- [1] Atlas. <https://www.bostondynamics.com/atlas>.
- [2] Fastest robot to solve a Rubik’s Cube. <https://www.guinnessworldrecords.com/world-records/fastest-robot-to-solve-a-rubiks-cube>.
- [3] <http://www.rcr3d.com/>. <https://mindcuber.com/>.
- [4] MindCuber. <https://mindcuber.com/>.
- [5] Partner Robot T-HR3. [https://www.toyota-global.com/innovation/partner\\_robot/robot/file/T-HR3\\_EN\\_0208.pdf](https://www.toyota-global.com/innovation/partner_robot/robot/file/T-HR3_EN_0208.pdf).
- [6] T. Abell and M. A. Erdmann. Stably supported rotations of a planar polygon with two frictionless contacts. In *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 1995, August 5 - 9, 1995, Pittsburgh, PA, USA*, pages 411–418, 1995.
- [7] Y. Aiyama, M. Inaba, and H. Inoue. Pivoting: A new method of grasplless manipulation of object by robot fingers. In *Proceedings of 1993 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 1993, Tokyo, Japan, July 26 - 30, 1993*, pages 136–143, 1993.
- [8] R. Antonova, S. Cruciani, C. Smith, and D. Kragic. Reinforcement learning for pivoting task. *CoRR*, abs/1703.00472, 2017.
- [9] D. Antotsiou, G. Garcia-Hernando, and T. Kim. Task-oriented hand motion retargeting for dexterous manipulation imitation. *CoRR*, abs/1810.01845, 2018.
- [10] T. Asfour, J. Schill, H. Peters, C. Klas, J. Bückner, C. Sander, S. Schulz, A. Kargov, T. Werner, and V. Bartenbach. Armar-4: A 63 dof torque controlled humanoid robot. In *2013 13th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, pages 390–396. IEEE, 2013.
- [11] Y. Bai and C. K. Liu. Dexterous manipulation using both palm and fingers. In *2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014*, pages 1560–1565, 2014.
- [12] G. Barth-Maron, M. W. Hoffman, D. Budden, W. Dabney, D. Horgan, D. TB, A. Muldal, N. Heess, and T. P. Lillicrap. Distributed distributional deterministic policy gradients. *CoRR*, abs/1804.08617, 2018.
- [13] A. Bicchi. Hands for dexterous manipulation and robust grasping: a difficult road toward simplicity. *IEEE Trans. Robotics and Automation*, 16(6):652–662, 2000.
- [14] A. Bicchi and R. Sorrentino. Dexterous manipulation through rolling. In *Proceedings of the 1995 International Conference on Robotics and Automation, Nagoya, Aichi, Japan, May 21-27, 1995*, pages 452–457, 1995.
- [15] M. Botvinick, S. Ritter, J. X. Wang, Z. Kurth-Nelson, C. Blundell, and D. Hassabis. Reinforcement learning, fast and slow. *Trends in cognitive sciences*, 2019.
- [16] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- [17] S. Carter, D. Ha, I. Johnson, and C. Olah. Experiments in handwriting with a neural network. *Distill*, 2016.
- [18] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox. Closing the Sim-to-Real Loop: Adapting Simulation Randomization with Real World Experience. *arXiv preprint 1810.05687*, 2018.
- [19] M. Cherif and K. K. Gupta. Planning quasi-static fingertip manipulations for reconfiguring objects. *IEEE Trans. Robotics and Automation*, 15(5):837–848, 1999.
- [20] M. Chociej, P. Welinder, and L. Weng. Orrb – openai remote rendering backend. *arXiv preprint 1906.11633*, 2019.
- [21] S. Christen, S. Stevsic, and O. Hilliges. Demonstration-guided deep reinforcement learning of control policies for dexterous human-robot interaction. *CoRR*, abs/1906.11695, 2019.
- [22] P. F. Christiano, Z. Shah, I. Mordatch, J. Schneider, T. Blackwell, J. Tobin, P. Abbeel, and W. Zaremba. Transfer from simulation to real world through learning deep inverse dynamics model. *CoRR*, abs/1610.03518, 2016.
- [23] I. Clavera, A. Nagabandi, R. S. Fearing, P. Abbeel, S. Levine, and C. Finn. Learning to adapt: Meta-learning for model-based control. *CoRR*, abs/1803.11347, 2018.
- [24] J. Clune. Ai-gas: Ai-generating algorithms, an alternate paradigm for producing general artificial intelligence. *arXiv preprint arXiv:1905.10985*, 2019.

- [25] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. AutoAugment: Learning Augmentation Policies from Data. *arXiv preprint 1805.09501*, may 2018.
- [26] A. Cully, J. Clune, D. Tarapore, and J.-B. Mouret. Robots that can adapt like animals. *Nature*, 521(7553):503, 2015.
- [27] W. Czarnecki, R. Pascanu, S. Osindero, S. M. Jayakumar, G. Swirszcz, and M. Jaderberg. Distilling policy distillation. *ArXiv*, abs/1902.02186, 2019.
- [28] N. C. Dafle and A. Rodriguez. Sampling-based planning of in-hand manipulation with external pushes. *CoRR*, abs/1707.00318, 2017.
- [29] N. C. Dafle, A. Rodriguez, R. Paolini, B. Tang, S. S. Srinivasa, M. A. Erdmann, M. T. Mason, I. Lundberg, H. Staab, and T. A. Fuhlbrigge. Extrinsic dexterity: In-hand manipulation with external forces. In *2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014*, pages 1578–1585, 2014.
- [30] Z. Doulgeri and L. Droukas. On rolling contact motion by robotic fingers via prescribed performance control. In *2013 IEEE International Conference on Robotics and Automation, Karlsruhe, Germany, May 6-10, 2013*, pages 3976–3981, 2013.
- [31] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel.  $RL^2$ : Fast reinforcement learning via slow reinforcement learning. *CoRR*, abs/1611.02779, 2016.
- [32] M. A. Erdmann. An exploration of nonprehensile two-palm manipulation. *I. J. Robotics Res.*, 17(5):485–503, 1998.
- [33] M. A. Erdmann and M. T. Mason. An exploration of sensorless manipulation. *IEEE J. Robotics and Automation*, 4(4):369–379, 1988.
- [34] P. Falco, A. Attawia, M. Saveriano, and D. Lee. On policy learning robust to irreversible events: An application to robotic in-hand manipulation. *IEEE Robotics and Automation Letters*, 3(3):1482–1489, 2018.
- [35] R. S. Fearing. Implementing a force strategy for object re-orientation. In *Proceedings of the 1986 IEEE International Conference on Robotics and Automation, San Francisco, California, USA, April 7-10, 1986*, pages 96–102, 1986.
- [36] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *CoRR*, abs/1703.03400, 2017.
- [37] A. Graves, M. G. Bellemare, J. Menick, R. Munos, and K. Kavukcuoglu. Automated curriculum learning for neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pages 1311–1320. JMLR.org, 2017.
- [38] S. Gu, E. Holly, T. P. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE International Conference on Robotics and Automation, ICRA 2017, Singapore, Singapore, May 29 - June 3, 2017*, pages 3389–3396, 2017.
- [39] M. Guo, A. Haque, D.-A. Huang, S. Yeung, and L. Fei-Fei. Dynamic task prioritization for multitask learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 270–287, 2018.
- [40] A. Gupta, C. Devin, Y. Liu, P. Abbeel, and S. Levine. Learning invariant feature spaces to transfer skills with reinforcement learning. *CoRR*, abs/1703.02949, 2017.
- [41] A. Gupta, B. Eysenbach, C. Finn, and S. Levine. Unsupervised meta-learning for reinforcement learning. *CoRR*, abs/1806.04640, 2018.
- [42] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine. Soft actor-critic algorithms and applications. *CoRR*, abs/1812.05905, 2018.
- [43] L. Han, Y. Guan, Z. X. Li, S. Qi, and J. C. Trinkle. Dextrous manipulation with rolling contacts. In *Proceedings of the 1997 IEEE International Conference on Robotics and Automation, Albuquerque, New Mexico, USA, April 20-25, 1997*, pages 992–997, 1997.
- [44] L. Han and J. C. Trinkle. Dextrous manipulation by rolling and finger gaiting. In *Proceedings of the IEEE International Conference on Robotics and Automation, ICRA-98, Leuven, Belgium, May 16-20, 1998*, pages 730–735, 1998.
- [45] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

- [46] R. Higo, Y. Yamakawa, T. Senoo, and M. Ishikawa. Rubik’s cube handling using a high-speed multi-fingered hand and a high-speed vision system. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6609–6614. IEEE, 2018.
- [47] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [48] W. H. Huang and M. T. Mason. Mechanics, planning, and control for tapping. *I. J. Robotics Res.*, 19(10):883–894, 2000.
- [49] J. Humplik, A. Galashov, L. Hasenclever, P. A. Ortega, Y. W. Teh, and N. Heess. Meta reinforcement learning as task inference. *CoRR*, abs/1905.06424, 2019.
- [50] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26):eaau5872, 2019.
- [51] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine. QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation. *ArXiv e-prints*, June 2018.
- [52] A. Kar, A. Prakash, M.-Y. Liu, E. Cameracci, J. Yuan, M. Rusiniak, D. Acuna, A. Torralba, and S. Fidler. Meta-Sim: Learning to Generate Synthetic Datasets. *arXiv preprint 1904.11621*, apr 2019.
- [53] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [54] V. Kumar, A. Gupta, E. Todorov, and S. Levine. Learning dexterous manipulation policies from experience and imitation. *CoRR*, abs/1611.05095, 2016.
- [55] V. Kumar, E. Todorov, and S. Levine. Optimal control with learned local models: Application to dexterous manipulation. In *2016 IEEE International Conference on Robotics and Automation, ICRA 2016, Stockholm, Sweden, May 16-21, 2016*, pages 378–383, 2016.
- [56] L. Lan, Z. Li, X. Guan, and P. Wang. Meta reinforcement learning with task embedding and shared policy. In *IJCAI*, 2019.
- [57] N. C. Landolfi, G. Thomas, and T. Ma. A model-based approach for sample-efficient multi-task reinforcement learning. *CoRR*, abs/1907.04964, 2019.
- [58] S. Levine and V. Koltun. Guided policy search. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pages 1–9, 2013.
- [59] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *I. J. Robotics Res.*, 37(4-5):421–436, 2018.
- [60] S. Levine, N. Wagener, and P. Abbeel. Learning contact-rich manipulation skills with guided policy search. In *IEEE International Conference on Robotics and Automation, ICRA 2015, Seattle, WA, USA, 26-30 May, 2015*, pages 156–163, 2015.
- [61] M. Li, Y. Bekiroglu, D. Kragic, and A. Billard. Learning of grasp adaptation through experience and tactile sensing. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems, Chicago, IL, USA, September 14-18, 2014*, pages 3339–3346, 2014.
- [62] M. Li, H. Yin, K. Tahara, and A. Billard. Learning object-level impedance control for robust grasping and dexterous manipulation. In *2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014*, pages 6784–6791, 2014.
- [63] Q. Li, M. Meier, R. Haschke, H. J. Ritter, and B. Bolder. Rotary object dexterous manipulation in hand: a feedback-based method. *IJMA*, 3(1):36–47, 2013.
- [64] T. Li, W. Xi, M. Fang, J. Xu, and M. Qing-Hu Meng. Learning to Solve a Rubik’s Cube with a Dexterous Hand. *arXiv e-prints*, page arXiv:1907.11388, Jul 2019.
- [65] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Thompson, S. Levine, and P. Sermanet. Learning latent plans from play. *CoRR*, abs/1903.01973, 2019.
- [66] R. R. Ma and A. M. Dollar. On dexterity and dexterous manipulation. In *15th International Conference on Advanced Robotics: New Boundaries for Robotics, ICAR 2011, Tallinn, Estonia, June 20-23, 2011.*, pages 1–7, 2011.
- [67] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. In *Robotics: Science and Systems XIII, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA, July 12-16, 2017*, 2017.
- [68] J. Mahler, M. Matl, X. Liu, A. Li, D. V. Gealy, and K. Goldberg. Dex-net 3.0: Computing robust robot suction grasp targets in point clouds using a new analytic model and deep learning. *CoRR*, abs/1709.06670, 2017.

- [69] T. Matiisen, A. Oliver, T. Cohen, and J. Schulman. Teacher-student curriculum learning. *arXiv preprint arXiv:1707.00183*, 2017.
- [70] B. Mehta, M. Diaz, F. Golemo, C. J. Pal, and L. Paull. Active Domain Randomization. *arXiv preprint 1904.04762*, 2019.
- [71] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. Meta-learning with temporal convolutions. *CoRR*, abs/1707.03141, 2017.
- [72] I. Mordatch, Z. Popovic, and E. Todorov. Contact-invariant optimization for hand manipulation. In *Proceedings of the 2012 Eurographics/ACM SIGGRAPH Symposium on Computer Animation, SCA 2012, Lausanne, Switzerland, 2012*, pages 137–144, 2012.
- [73] A. Nagabandi, K. Konoglie, S. Levine, and V. Kumar. Deep Dynamics Models for Learning Dexterous Manipulation. In *Conference on Robot Learning (CoRL)*, 2019.
- [74] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [75] Nordic Semiconductor. nRF52832 Product Specification v1.1. Technical report, Nordic Semiconductor, 2016.
- [76] A. M. Okamura, N. Smaby, and M. R. Cutkosky. An overview of dexterous manipulation. In *Proceedings of the 2000 IEEE International Conference on Robotics and Automation, ICRA 2000, April 24-28, 2000, San Francisco, CA, USA*, pages 255–262, 2000.
- [77] C. Olah, A. Satyanarayan, I. Johnson, S. Carter, L. Schubert, K. Ye, and A. Mordvintsev. The building blocks of interpretability. *Distill*, 2018. <https://distill.pub/2018/building-blocks>.
- [78] OpenAI. OpenAI Five. <https://blog.openai.com/openai-five/>, 2018.
- [79] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba. Learning dexterous in-hand manipulation. *CoRR*, 2018.
- [80] E. Parisotto, S. Ghosh, S. B. Yalamanchi, V. Chinnabireddy, Y. Wu, and R. Salakhutdinov. Concurrent meta reinforcement learning. *CoRR*, abs/1903.02710, 2019.
- [81] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. *CoRR*, abs/1710.06537, 2017.
- [82] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel. Asymmetric actor critic for image-based robot learning. *arXiv preprint arXiv:1710.06542*, 2017.
- [83] L. Pinto, J. Davidson, and A. Gupta. Supervision via competition: Robot adversaries for learning tasks. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1601–1608. IEEE, 2017.
- [84] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta. Robust adversarial reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2817–2826. JMLR. org, 2017.
- [85] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin, M. Chociej, P. Welinder, et al. Multi-goal reinforcement learning: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018.
- [86] A. Prakash, S. Bochoon, M. Brophy, D. Acuna, E. Cameracci, G. State, O. Shapira, and S. Birchfield. Structured Domain Randomization: Bridging the Reality Gap by Context-Aware Synthetic Data. *arXiv preprint 1810.10093*, oct 2018.
- [87] A. Rajeswaran, V. Kumar, A. Gupta, J. Schulman, E. Todorov, and S. Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *CoRR*, abs/1709.10087, 2017.
- [88] T. Rokicki and M. Davidson. God’s Number is 26 in the Quarter-Turn Metric. <https://cube20.org/qtm/>, 2014.
- [89] S. Ross, G. J. Gordon, and J. A. Bagnell. No-regret reductions for imitation learning and structured prediction. *CoRR*, abs/1011.0686, 2010.
- [90] N. Ruiz, S. Schuler, and M. Chandraker. Learning To Simulate. *arXiv preprint 1810.02513*, 1810.02513:1–12, 2018.
- [91] D. Rus. Dexterous rotations of polyhedra. In *Proceedings of the 1992 IEEE International Conference on Robotics and Automation, Nice, France, May 12-14, 1992*, pages 2758–2763, 1992.
- [92] D. Rus. In-hand dexterous manipulation of piecewise-smooth 3-d objects. *I. J. Robotics Res.*, 18(4):355–381, 1999.



- [93] A. A. Rusu, M. Vecerik, T. Rothörl, N. Heess, R. Pascanu, and R. Hadsell. Sim-to-real robot learning from pixels with progressive nets. In *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, pages 262–270, 2017.
- [94] F. Sadeghi and S. Levine. CAD2RL: real single-image flight without a single real image. In *Robotics: Science and Systems XIII, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA, July 12-16, 2017*, 2017.
- [95] S. Sæmundsson, K. Hofmann, and M. P. Deisenroth. Meta reinforcement learning with latent variable gaussian processes. In *UAI*, 2018.
- [96] Y. Sakagami, R. Watanabe, C. Aoyama, S. Matsunaga, N. Higaki, and K. Fujimura. The intelligent asimo: System overview and integration. In *IEEE/RSJ international conference on intelligent robots and systems*, volume 3, pages 2478–2483. IEEE, 2002.
- [97] N. Sawasaki and H. INOUE. Tumbling objects using a multi-fingered robot. *Journal of the Robotics Society of Japan*, 9(5):560–571, 1991.
- [98] J. Schmidhuber. POWERPLAY: training an increasingly general problem solver by continually searching for the simplest still unsolvable problem. *CoRR*, abs/1112.5309, 2011.
- [99] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [100] ShadowRobot. ShadowRobot Dexterous Hand. <https://www.shadowrobot.com/products/dexterous-hand/>, 2005.
- [101] J. Shi, J. Z. Woodruff, P. B. Umbanhowar, and K. M. Lynch. Dynamic in-hand sliding manipulation. *IEEE Trans. Robotics*, 33(4):778–795, 2017.
- [102] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–503, 2016.
- [103] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. P. Lillicrap, K. Simonyan, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *CoRR*, abs/1712.01815, 2017.
- [104] R. K. Srivastava, B. R. Steunebrink, and J. Schmidhuber. First experiments with powerplay. *CoRR*, abs/1210.8385, 2012.
- [105] K. Tahara, S. Arimoto, and M. Yoshida. Dynamic object manipulation using a virtual frame by a triple soft-fingered robotic hand. In *IEEE International Conference on Robotics and Automation, ICRA 2010, Anchorage, Alaska, USA, 3-7 May 2010*, pages 4322–4327, 2010.
- [106] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. *CoRR*, abs/1804.10332, 2018.
- [107] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. *arXiv preprint arXiv:1703.06907*, 2017.
- [108] J. Tobin, W. Zaremba, and P. Abbeel. Domain randomization and generative models for robotic grasping. *CoRR*, abs/1710.06425, 2017.
- [109] E. Todorov, T. Erez, and Y. 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.
- [110] P. Tournassoud, T. Lozano-Pérez, and E. Mazer. Regrasping. In *Proceedings of the 1987 IEEE International Conference on Robotics and Automation, Raleigh, North Carolina, USA, March 31 - April 3, 1987*, pages 1924–1928, 1987.
- [111] M. Tsou. Kociemba. <https://github.com/muodov/kociemba>.
- [112] E. Tzeng, C. Devin, J. Hoffman, C. Finn, X. Peng, S. Levine, K. Saenko, and T. Darrell. Towards adapting deep visuomotor representations from simulated to real environments. *CoRR*, abs/1511.07111, 2015.
- [113] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [114] H. van Hoof, T. Hermans, G. Neumann, and J. Peters. Learning robot in-hand manipulation with tactile features. In *15th IEEE-RAS International Conference on Humanoid Robots, Humanoids 2015, Seoul, South Korea, November 3-5, 2015*, pages 121–127, 2015.

- [115] Q. Vuong, S. Vikram, H. Su, S. Gao, and H. I. Christensen. How to pick the domain randomization parameters for sim-to-real transfer of reinforcement learning policies? *arXiv preprint 1903.11774*, pages 1–3, 2019.
- [116] J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Kumaran, and M. Botvinick. Learning to reinforcement learn. *CoRR*, abs/1611.05763, 2016.
- [117] R. Wang, J. Lehman, J. Clune, and K. O. Stanley. Poet: open-ended coevolution of environments and their optimized solutions. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 142–151. ACM, 2019.
- [118] Y. You, I. Gitman, and B. Ginsburg. Scaling SGD batch size to 32k for imagenet training. *CoRR*, abs/1708.03888, 2017.
- [119] W. Yu, J. Tan, C. K. Liu, and G. Turk. Preparing for the unknown: Learning a universal policy with online system identification. In *Robotics: Science and Systems XIII, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA, July 12-16, 2017*, 2017.
- [120] S. Zakharov, W. Kehl, and S. Ilic. DeceptionNet: Network-Driven Domain Randomization. *arXiv preprint 1904.02750*, apr 2019.
- [121] M. Zhang, S. Vikram, L. Smith, P. Abbeel, M. J. Johnson, and S. Levine. SOLAR: deep structured latent representations for model-based reinforcement learning. *CoRR*, abs/1808.09105, 2018.
- [122] H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar. Dexterous manipulation with deep reinforcement learning: Efficient, general, and low-cost. *CoRR*, abs/1810.06045, 2018.
- [123] Y. Zhu, Z. Wang, J. Merel, A. A. Rusu, T. Erez, S. Cabi, S. Tunyasuvunakool, J. Kramár, R. Hadsell, N. de Freitas, and N. Heess. Reinforcement and imitation learning for diverse visuomotor skills. *CoRR*, abs/1802.09564, 2018.

# Appendices

## A System Monitoring

Reliability of the physical setup was one of the key challenges mentioned in [79]. To address this issue we built new monitoring system which ensures experiments were performed on a fully functional robot and that physical setup was consistent for all experiments. New system includes recording, visualization, alerting and persistence of all operations performed on the physical robots. New debugging and investigation capabilities were unlocked and allowed us to detect robot breakage in real time (e.g. breaking tendons), find regressions in the physical setup and identify differences between policies deployed on the physical setup and in the simulation.

We used InfluxDB<sup>21</sup> database to do real-time tracking and persistence of native resolution sensor data and detailed information about the context in which experiments were conducted.

We used Grafana<sup>22</sup>, with InfluxDB as a data source, to display information on policy performance on the task, to visualize data from robot sensors and to show the health status of robotic components.

## B Randomizations

Detailed listings of the randomizations used in policy and vision training are given in Tables 9 and 11. Below, we describe how the randomizations are parameterized in ADR.

A simple randomization for environment parameter  $\lambda_i$ , such as gravity, is parameterized by two ADR parameters corresponding to the boundaries of its randomization range:  $\phi_{L,i}$  and  $\phi_{H,i}$ . Most randomizations that we used were simple. A more complex randomization such as observation noise is parameterized by several ADR parameters.

In the following, we denote default (non-randomized) values by using the subscript  $(\cdot)_0$ . We use  $\lambda_i, \lambda_j, \dots$  to denote randomized environment parameters, which are parameterized in ADR by boundary values  $\phi_{L,i}, \phi_{H,i}, \phi_{L,j}, \phi_{H,j}$ , etc. Lastly, define  $g(x) \triangleq \exp(x - 1.0)$ .

### B.1 Simulator Physics Randomizations

Simulator physics parameters are randomized by using either *generic* or *custom* randomizers.

#### B.1.1 Generic randomizers

A generic randomizer is specified by a *noise mode* and a scaling factor  $\alpha$ . The noise modes used in this work are listed in Table 10. The variable  $x$  denotes a simulator physics parameter being randomized.

#### B.1.2 Custom randomizers

We also used several custom randomizers that perform slightly different randomization methods than generic randomizers. Again,  $x$  denotes a simulator physics parameter being randomized.

**Cube and robot friction.** The slide, spin, or roll friction of a cube or robot body is randomized by

$$x = x_0 e^{w\lambda_i},$$

where  $w$  is a fixed weight specific for the type of friction. The weights for robot friction are 1.0 for all types. The weights for cube friction is  $w = 1.0$  for slide and  $w = 2.0$  for spin and roll.

**Cube size.** Cube size is randomized by

$$x = x_0 e^{0.15\lambda_i}$$

<sup>21</sup>InfluxDB is a time series database designed to handle high write and query loads. See <https://www.influxdata.com/products/influxdb-overview> for more information.

<sup>22</sup>Grafana is an open-source, general purpose dashboard and graph composer, which runs as a web application. See <https://grafana.com> for more information.

Table 9: Randomizations used to train manipulation policies with ADR. For simulator physics randomizations with generic randomizers, values in the parenthesis denote (noise mode,  $\alpha$ ). Generic randomizers without parenthesis used default values (MG, 1.0).

Category	All policies	Reorientation policy	Rubik’s cube policy
Simulator physics (generic)	Actuator force range Actuator gain prm Body inertia Geom size robot spatial Tendon length spring (M, 0.75) Tendon stiffness (M, 0.75)	Dof armature Dof damping Dof friction loss Geom friction Geom gap (M, 0.03)	Body position (AG, 0.02) Dof armature cube Dof armature robot Dof damping cube Dof damping robot Dof friction loss cube Dof friction loss robot Geom gap cube (AU, 0.01) Geom gap robot (AU, 0.01) Geom pos cube (AG, 0.002) Geom pos robot (AG, 0.002) Geom margin cube (AG, 0.0005) Geom margin robot (AG, 0.0005) Geom solimp (M, 1.0) Geom solref (M, 1.0) Joint stiffness robot (UAG, 0.005)
Simulator physics (custom)	Body mass Cube size Tendon range		Friction robot Friction cube
Custom physics	Action latency Backlash Joint margin Joint range Time step		Action noise Time step variance
Adversary	Adversary		
Observation			Observation

**Joint and tendon limits.** Both the lower and upper limits on joint and tendon ranges are randomized by

$$x = x_0 + n \text{ for } n \sim \mathcal{N}(0, 0.1g(|\lambda_i|)).$$

## B.2 Custom Physics Randomizations

Custom physics models were used to capture physical effects that are not modelled by the simulator. These models range from simple time delays to motor backlash.

**Action delay.** Action delay  $d$  in milliseconds is modelled and randomized by

$$d = |\lambda_i|n_0n_1 \text{ for } N_0 \sim \mathcal{N}(1, |\lambda_j|^2), N_1 \sim \mathcal{N}(1, |\lambda_k|^2),$$

where  $n_0$  is sampled once per episode,  $n_1$  is sampled once per step.

**Action latency.** Action latency  $l$  in time steps is randomized by

$$l \sim C[\lambda_i],$$

where  $C[n]$  is the uniform categorical distribution over  $n$  elements.

Table 10: Generic randomizer noise modes and their abbreviations.  $x$  denotes a simulator parameter being randomized and  $x_0$  its default value.

Noise Mode	Abbreviation	Expression
Additive Gaussian	AG	$x_0 +  N $ for $N \sim \mathcal{N}(g(\alpha\lambda_i), g( \alpha\lambda_j )^2)$
Unbiased Additive Gaussian	UAG	$x_0 + N$ for $N \sim \mathcal{N}(0, g( \alpha\lambda_i )^2)$
Multiplicative	M	$x_0 e^N$ for $N \sim \mathcal{N}(\alpha\lambda_i,  \alpha\lambda_j ^2)$

**Action noise.** Action noise  $a$  is randomized by

$$a = a_0 n_0 + n_1 + n_2 \text{ for } n_0 \sim \mathcal{N}(1, g(|\lambda_i|)^2), n_1 \sim \mathcal{N}(0, g(|\lambda_j|)^2), n_2 \sim \mathcal{N}(0, g(|\lambda_k|)^2),$$

where  $n_0, n_1$  are sampled once per episode,  $n_2$  is sampled once per step.

**Motor Backlash.** As described in [79, C.2] our motor backlash implementation contains two parameters  $\delta_{\pm 1}$ . They are randomized by

$$\delta_{\pm 1} = e^{n\delta_{\pm 1,0}} \text{ for } n \sim \mathcal{N}(1, \lambda_i^2).$$

**Gravity.** Gravity vector  $g$  is randomized by

$$g = g_0 + ug(\lambda_i)$$

for  $u \in \mathbb{R}^3$  a random vector uniformly distributed over the unit sphere.

**Joint margin.** Joint margin  $m$  is randomized by

$$m = m_0 n \text{ for } n \sim U[0, 0.15g(\lambda_i)].$$

**Time step.** Simulation time step  $t$  is randomized by

$$t = e^{0.6\lambda_i}(t_0 + ne^{\lambda_j}) \text{ for } n \sim \text{Exp}(1/\kappa), \kappa \sim U[1250, 10000].$$

The *time step variance* randomization is equivalent to the time step randomization with  $\lambda_i = 0$ .

### B.3 Random Network Adversary (RNA)

ADR allows us to automatically choose randomization ranges and allows us to train on an ever-growing distribution over environments. However, one problem that ADR faces is that there might be unmodeled effects in the target domain. Since they are not modeled, ADR cannot randomize them.

To resolve this, we use an adversarial approach similar to [83, 84]. However, instead of training the adversary in a zero-sum fashion, we use networks with randomly sampled weights. We re-sample these weights at the beginning of every episode and keep them fixed throughout the episode. Our adversarial approach has two different means of perturbing the policy:

- *Action perturbations.* Let  $\pi_{\text{adv}}$  denote an adversarial policy with the same action space as the policy controlling the robot (usually denoted just as  $\pi$  but denoted as  $\pi_{\text{robot}}$  in this section for clarity). We then use the following simple convex combination to obtain the action we execute in simulation:  $\mathbf{a}_t = (1 - \alpha)\mathbf{a}_{\text{robot}} + \alpha\mathbf{a}_{\text{adv}}$ , where  $\alpha \in [0, 1]$ ,  $\mathbf{a}_{\text{robot}} \sim \pi_{\text{robot}}(\mathbf{o}_t)$  and  $\mathbf{a}_{\text{adv}} \sim \pi_{\text{adv}}(\mathbf{o}_t)$ .  $\alpha$  controls the influence of the adversary, with  $\alpha = 0$  meaning no adversarial perturbation and  $\alpha = 1$  meaning only adversarial perturbations. Given a sufficiently powerful adversarial policy, action perturbations allow us to model arbitrary effects in the actuation system of the robot.
- *Force/torque perturbations.* Similar to action perturbations, we apply force/torque perturbations to all bodies of the object being manipulated.<sup>23</sup> More concretely, we generate a  $6D$  vector for each body where the first 3 dimensions define a force vector and the last 3 dimensions define a torque vector for the same body. We scale the force vector by the body’s mass and the torques by the body’s corresponding inertia (we use a diagonal inertia matrix and scale each element of the torque vector by the corresponding diagonal entry) to normalize the effect each force/torque has. We then scale all force/torque perturbations by a scalar  $\beta \in \mathbb{R}_{\geq 0}$ , thus allowing us to control how much influence force/torque perturbations have.

<sup>23</sup>For the block reorientation task this is just a single body but for the Rubik’s cube we apply perturbations to each cublet.

Given this formulation, we can use ADR to automatically adjust  $\alpha$  and  $\beta$  over time without the need for hand-tuning (which would be exceedingly difficult). Furthermore, we can sample  $\alpha$  and  $\beta$  iid per coordinate (i.e. per joint and per body, respectively), this having simulations in which some joints or bodies are “more adversarial” than others.

However, one question still remains: How do we obtain the perturbations, i.e. how do we implement the adversarial policy? We experimented with sampling random perturbations, using a zero-sum self-play setup where we train an adversarial policy directly using either feed-forward or recurrent networks. Surprisingly we found that a very simple method outperforms them in terms of sim2sim transfer: Random networks.

More concretely, we use a simple feed-forward network with 3 hidden layers of 265 units each. After each hidden layer we use a ReLU non-linearity. We use a discretized action space with 31 bins per dimensions. We do this to ensure that randomly initializing the network leads uniform probability of selecting one of the 31 bins (if we would use a continuous space and a tanh non-linearity, actions would always be close to 0). We re-initialize the adversarial network at the beginning of each episode and keep it fixed after. The inputs of the adversarial network are the noisy fingertip positions as well as cube position, orientation, and face angle configuration. This idea is also related to random network distillation, where a randomly initialized CNN is used to help with exploration [16].

We found that this approach works very well, outperforming more sophisticated methods like adversarial training. We believe this is because diversity is what ultimately aids in transfer. Even though a trained adversary is a stronger opponent, random networks provide the maximum amount of diversity.

#### B.4 Observation Randomizations

We add both correlated and uncorrelated noise to observations. Correlated noise is sampled once per episode and uncorrelated noise is sampled once per step. Given default noise standard deviations  $a_0, b_0, c_0$ , fixed for each observation element  $o$ , we randomize each observation by

$$o = o_0 n_0 + n_1 + n_2 \text{ for } n_0 \sim \mathcal{N}(1, (a_0 e^{\lambda_i})^2), n_1 \sim \mathcal{N}(0, (b_0 e^{\lambda_i})^2), n_2 \sim \mathcal{N}(0, (c_0 e^{\lambda_j})^2),$$

where  $n_0, n_1$  are sampled once per episode,  $n_2$  is sampled once per step.

#### B.5 Visual Randomizations

There are two categories of randomizations for vision training, (a) variations in the scenes that are rendered by ORRB [20] and (b) TensorFlow distortion operations applied on the images. ADR controls to what extent each category of randomizations affect the appearance of final images that are fed into the model with different ADR parameters. See the full list in Table 11.

<sup>24</sup><https://docs.unity3d.com/Packages/com.unity.postprocessing@2.1/manual/Bloom.html>

<sup>25</sup><https://docs.unity3d.com/Packages/com.unity.postprocessing@2.1/manual/Ambient-Occlusion.html>

Table 11: Parameter randomizations used in vision model ADR training.

Category	Randomizer	Parameter
ORRB [20]	Camera Randomizer	Position perturbation distance Rotation perturbation magnitude Field of view (fov) perturbation
	Light Randomizer	Individual light intensity Total intensity of lights in the scene Spotlight angle
	Lighting Rig	Number of lights Light distance Light height
	Material Fixed Hue Randomizer	Hue perturbation radius Saturation perturbation radius Value perturbation radius Range of saturation Range of value Range of glossiness Range of metallic
	Post Processing Randomizer	Hue shift of all colors Saturation of all colors Contrast of all colors Brightness of all colors Color temperature to set the white balance to Range of tint Strength of bloom filter <sup>24</sup> Extent of veiling effects Degree of darkness added by ambient occlusion <sup>25</sup> Grain strength Grain article size
TensorFlow	—	<code>tf.image.adjust_hue</code> <code>tf.image.adjust_saturation</code> <code>tf.image.adjust_brightness</code> <code>tf.image.adjust_contrast</code> Add random Gaussian noise Add random Gaussian noise that is same across channels

## C Implementation Details

### C.1 Goal Generation

Algorithm listings for goal generations for both tasks are given here. It is important to note that in both cases random orientation is sampled in such way that one face of the block or cube is pointing directly upwards. For reference, as was described earlier, the cube or block orientations are considered aligned if any face is pointing upwards within a tolerance of 0.4 radians. Cube faces are considered aligned if all six are within 0.1 radians from a straight angle configuration.



**Algorithm 3** Block reorientation goal generation

---

$goal\_orientation \leftarrow \text{RANDOMUNIFORMUPWARDORIENTATION}$  ▷ Sample random quaternion  
 $\text{SETORIENTATIONGOAL}(goal\_orientation)$

---

**Algorithm 4** Rubik’s cube goal generation

---

**Require:**  $cube\_quat$  ▷ Quaternion describing current cube orientation  
**Require:**  $face\_angles$  ▷ Current rotation angles of six faces of the cube  
 $is\_aligned \leftarrow \text{ISORIENTATIONALIGNED}(cube\_quat)$  **and**  $\text{ISFACEALIGNED}(face\_angles)$  ▷ Check if cube is in the *aligned* state  
 $u \leftarrow U(0, 1)$  ▷ Sample random number uniformly between 0 and 1  
  
**if**  $is\_aligned$  **and**  $(u < 0.5)$  **then** ▷ If cube is aligned then with 50% probability we perform face rotation  
     $side \leftarrow \text{RANDOM}(\{CW, CCW\})$  ▷ Choose clockwise or counterclockwise  
     $goal\_face\_angles \leftarrow \text{ROTATETOPFACE}(side, face\_angles)$  ▷ Rotate top face by 90 degrees  
     $\text{SETFACEGOAL}(goal\_face\_angles)$   
     $\text{SETORIENTATIONGOAL}(NULL)$   
**else** ▷ Otherwise, perform cube flip  
     $aligned\_face\_angles \leftarrow \text{ALIGN}(face\_angles)$  ▷ Align face angles to closest straight angle configuration  
     $goal\_orientation \leftarrow \text{RANDOMUNIFORMUPWARDORIENTATION}$  ▷ Sample random quaternion  
     $\text{SETFACEGOAL}(aligned\_face\_angles)$   
     $\text{SETORIENTATIONGOAL}(goal\_orientation)$   
**end if**

---

**C.2 Neural Network Inputs**

Table 12 lists the policy and value function inputs for the block reorientation task, respectively. The number of parameters for the Rubik’s cube task network is given in Table 13.

Table 12: Observations for the block reorientation task of the policy and value networks, respectively.

Input	Dimensionality	Policy network	Value network
Fingertip positions	15D	×	✓
Noisy fingertip positions	15D	✓	✓
Cube position	3D	×	✓
Noisy block position	3D	✓	✓
Block orientation	4D (quaternion)	×	✓
Noisy block orientation	4D (quaternion)	✓	✓
Goal orientation	4D (quaternion)	✓	✓
Relative goal orientation	4D (quaternion)	×	✓
Noisy relative goal orientation	4D (quaternion)	✓	✓
Hand joint angles	48D <sup>26</sup>	×	✓
All simulation positions & orientations (qpos)	38D	×	✓
All simulation velocities (qvel)	36D	×	✓

---

<sup>26</sup>Angles are encoded as sin and cos, i.e. this doubles the dimensionality of the underlying angle.

Table 13: Network size measured by the number of parameters for the Rubik’s cube network (parameters for the block reorientation task are the same except for the input embeddings).

Type	Number of parameters
Nontrainable parameters	1539
Policy network	13,863,132
Value function network	13,638,657
Input embeddings	267,776

### C.3 Hyperparameters

Hyperparameters used in our PPO policy are listed in Table 14.

Table 15 shows hyperparameters used in ADR for policy training.

Vision training hyperparameters are given in Table 16 and the details of the model architecture are given in Table 17. The target errors for our loss weight auto-balancer and vision ADR are listed in Table 18. The increase and decrease thresholds in vision ADR are 70% and 30% respectively.

We also apply data augmentation for vision training. More specifically, we leave the object pose as is with 20% probability, rotate the object by  $90^\circ$  around its main axes with 40% probability, and “jitter” the object by adding Gaussian noise to both the position and rotation independently with 40% probability.

Table 14: Hyperparameters used for PPO.

Hyperparameter	Value
hardware configuration - block	32 NVIDIA V100 GPUs + 12’800 CPU cores
hardware configuration - Rubik’s Cube	64 NVIDIA V100 GPUs + 29’440 CPU cores
action distribution	categorical with 11 bins for each one of 20 action coordinates
discount factor $\gamma$	0.998
Generalized Advantage Estimation $\lambda$	0.95
entropy regularization coefficient	varying 0.01 – 0.0025
PPO clipping parameter $\epsilon$	0.2
optimizer	Adam [53]
learning rate	varying $3 \times 10^{-4}$ – $1 \times 10^{-4}$
batch size (per GPU)	5120 chunks x 10 transitions = 51’200 frames
Sample reuse (experience replay)	3
Value loss weight	1.0
L2 regularization weight	$10^{-6}$

Table 15: Hyperparameters used for ADR.

Hyperparameter	Value
Maximum value of a $\phi$	4.0
Boundary sampling probability	0.5
ADR step size $\Delta$	0.02
ADR increase threshold $t_H$	20
ADR decrease threshold $t_L$	10
Performance queue length $m$	240

Table 16: Hyperparameters used for the vision model training.

Hyperparameter	Value
optimization hardware	8 NVIDIA V100 GPUs + 64 vCPU cores
rendering hardware	16 $\times$ 1 NVIDIA V100 GPU + 8 vCPU core machines
optimizer	LARS [118]
learning rate	0.5, halved every 40 000 batches
minibatch size	$256 \times 3 = 768$ RGB images
image size	$200 \times 200$ pixels
weight decay regularization	0.0001
number of training batches	300 000
network architecture	shown in Figure 14

Table 17: Hyperparameters for the vision model architecture.

Layer	Details
Input RGB Image	$200 \times 200 \times 3$
Batch Norm	
Conv2D	64 filters, $5 \times 5$ , stride 2, no padding
Dropout	keep probability = 0.9
Max Pooling	$2 \times 2$ , stride 1, no padding
Batch Norm	
ResNet [45] Bottleneck	3 blocks, 64 filters, $1 \times 1$ , stride 1
Dropout	keep probability = 0.9
ResNet [45] Bottleneck	4 blocks, 128 filters, $2 \times 2$ , stride 2
ResNet [45] Bottleneck	6 blocks, 256 filters, $2 \times 2$ , stride 2
ResNet [45] Bottleneck	3 blocks, 512 filters, $2 \times 2$ , stride 2
Max Pooling	$2 \times 2$ , stride 1, no padding
Flatten	
Concatenate	all 3 image towers combined
Dropout	keep probability = 0.9
Fully Connected	512 units
Batch Norm	
Fully Connected	256 units
Batch Norm	
Fully Connected	output, 3 pos. + 4 quat. + 180 top face angle + 90 active angles + 3 active axis = 280 units

Table 18: Target error thresholds for different prediction labels in the loss weight auto-balancer and ADR of the vision model. The active axis target error is set to 0.5 because the error is binary for each individual sample.

Label	Target Errors	
	Auto Balancer	Vision ADR
Orientation	$3^\circ$	$5^\circ$
Position	3 mm	5 mm
Top face angle	$3^\circ$	$5^\circ$
Active axis	0.5	0.5
Active face angles	$3^\circ$	$5^\circ$

## C.4 Vision Post-Processing

Algorithm 5 describes the post-processing logic that we apply when using the vision model for face angle predictions.

---

**Algorithm 5** Vision model post-processing logic for tracking the full state of a Rubik’s cube on the physical robot.

---

**Require:**  $base\_angles[0..5]$  ▷ 6 angles in total  
**Require:**  $tracked\_angles[0..5]$   
**Require:**  $cube\_quat$  ▷ Predicted cube orientation  
**Require:**  $active\_axis$  ▷ Predicted active axis,  $\in (0, 1, 2)$   
**Require:**  $active\_angles[0, 1]$  ▷ Predicted two active face angles in  $[-\pi/4, \pi/4]$   
**Require:**  $top\_angle$  ▷ Predicted top face angle in  $[-\pi, \pi]$   
 $k \leftarrow \text{EXTRACTTOPFACEID}(cube\_quat)$   
 $is\_aligned \leftarrow \text{ISORIENTATIONALIGNED}(cube\_quat)$   
**for**  $i := 0$  to  $5$  **do**  
  **if**  $i = k$  **and**  $is\_aligned$  **then** ▷ If this face is placed on top  
     $base\_angles[i] \leftarrow top\_angle$   
     $tracked\_angles[i] \leftarrow top\_angle$   
  **else**  
     $a, b \leftarrow \text{divmod}(i, 2)$   
    **if**  $a = active\_axis$  **then** ▷ If this is a active face  
       $tracked\_angles[i] \leftarrow \text{MOVEANGLE}(base\_angles[i], active\_angles[b])$  ▷ See Algorithm 6  
    **end if**  
  **end if**  
**end for**

---



---

**Algorithm 6** MOVEANGLE: Move the base face angle to the target by a modulo  $\pi/2$  delta angle.

---

**Require:**  $a$  ▷ Base face angle  
**Require:**  $b$  ▷ Target face angle  
 $\delta \leftarrow \text{mod}(b - a + \pi/4, \pi/2) - \pi/4$   
**return**  $a + \delta$

---

## D Full Results

### D.1 Simulation Calibration

We test how much of an impact simulation calibration has. We use a simplified version of the Rubik’s cube: The face cube. In this version, the Rubik’s cube is constrained to only have 2 degrees of freedom with the other degrees disabled. We evaluate a policy trained on the old simulation and on the new simulation (i.e. with coupling and dynamics calibration). We calibrate the following MuJoCo parameters that are related to joint movements: `dof_damping`, `jnt_range`, `geom_size`, `tendon_range`, `tendon_lengths`, `spring`, `tendon_stiffness`, `actuator_forcerange`, and `actuator_gainprm`. Results are in Table 19.

Table 19: Performance of two policies trained on two different environments on a simplified version of the Rubik’s cube task. We run 10 trials for each policy on the physical setup and report the mean and median number of face rotations and mean number of policy steps per face rotation.

Policy	Successes		Steps Per Success
	Mean	Median	Mean
Original Simulation	$4.8 \pm 2.89$	5.00	321.10
New Simulation (Coupling Model + Calibration)	<b><math>14.30 \pm 6.56</math></b>	<b>15.50</b>	<b>252.50</b>

### D.2 Physical Trials

The full statistics of our physical trials are available in Table 20 and Table 21.

Table 20: Performance of different policies on the block reorientation task. We evaluate each policy on the real robot for 10 trials. For each individual trial, we report the mean  $\pm$  standard error and median number of successes, the time per success in second, as well as the full list of success counts.

Policy	Mean	Median	Time per Success	Trials (sorted)
Baseline (from [79])	$18.8 \pm 5.4$	13.0	—	50, 41, 29, 27, 14, 12, 6, 4, 4, 1
Baseline (re-run of [79])	$4.0 \pm 1.7$	2.0	17.24 sec	17, 10, 3, 3, 2, 2, 1, 1, 1, 0
Manual DR	$2.7 \pm 1.1$	1.0	31.65 sec	11, 5, 4, 3, 1, 1, 1, 1, 0, 0
ADR (Small)	$1.4 \pm 0.9$	0.5	67.49 sec	9, 2, 1, 1, 1, 0, 0, 0, 0, 0
ADR (Medium)	$3.2 \pm 1.2$	2.0	29.50 sec	12, 7, 4, 3, 2, 2, 1, 1, 0, 0
ADR (Large)	$13.3 \pm 3.6$	11.5	4.45 sec	38, 25, 17, 16, 12, 11, 5, 4, 3, 2
ADR (XL)	$16.0 \pm 4.0$	12.5	<b>5.10</b> sec	42, 28, 27, 16, 13, 12, 8, 7, 5, 2
ADR (XXL)	<b><math>32.0 \pm 6.4</math></b>	<b>42.0</b>	5.18 sec	<b>50, 50, 50, 50</b> , 43, 41, 13, 12, 6, 5

Table 21: Performance of different policies on the Rubik’s cube for a fixed fair scramble goal sequence. We evaluate each policy on the real robot for 10 trials. For each individual trial, we report the mean  $\pm$  standard error and median number of successes, the time per success in second, as well as the full list of success counts.

Policy	Sensing	Mean	Median	Time per Success (Flip) (Rotation)		Trials (sorted)
Manual DR	Giiker + Vision	$1.8 \pm 0.4$	2.0	102.92 sec	35.60 sec	4, 3, 2, 2, 2, 2, 1, 0, 0
ADR	Giiker + Vision	$3.8 \pm 1.0$	3.0	66.97 sec	30.65 sec	8, 8, 8, 5, 4, 2, 2, 1, 0, 0
ADR (XL)	Giiker + Vision	$17.8 \pm 4.2$	12.5	<b>6.44</b> sec	<b>10.98</b> sec	44, 38, 24, 17, 14, 11, 9, 8, 7, 6
ADR (XXL)	Giiker + Vision	<b><math>26.8 \pm 4.9</math></b>	<b>22.0</b>	6.55 sec	11.79 sec	<b>50, 50</b> , 42, 24, 22, 22, 21, 19, 13, 5
ADR (XXL)	Vision	$12.8 \pm 3.4$	10.5	9.71 sec	13.23 sec	31, 25, 21, 18, 17, 4, 3, 3, 3, 3

### D.3 Vision Model Performance

Our vision model for estimating the full state of a Rubik’s cube outputs five labels. The prediction errors are evaluated on both simulated and real data. Orientation error is computed as rotational distance over a quaternion representation. Position error is the euclidean distance in 3D space, in millimeters. Face angle error is measured in degree ( $^{\circ}$ ). For active axis, we measure the percentage of incorrect predicted label (equivalent to  $1.0 - \text{accuracy}$ ). Note that the high error on the active axis prediction is mostly due to the fact that the cube is often aligned with very tiny active face angles and hence the active axis is less pronounced.

- **Orientation:** Cube orientation in quaternion.
- **Position:** Cube position in millimetre.
- **Top angle:** The full angle of the top face in  $[-\pi, \pi]$ .
- **Active axis:** Which one out of three axes is active (i.e., has non-aligned faces).
- **Active angles:** The angles of two faces relevant for the active axis in  $[-\pi/4, \pi/4]$ .

### D.4 Meta-Learning

The meta-learning results for face rotations are available in Figure 21. The results are comparable to the cube flips presented in the main body of the paper. We include them in the appendix for completeness.

Table 22: Performance of vision models in all our experiments for the Rubik’s cube prediction task, including experiment in the ablation study and models trained at different ADR entropy levels.

	Errors	Baseline	No DR	No focal loss	Non-discrete angles	ADR (S)	ADR (M)	ADR (L)
Sim	<b>Orientation</b>	6.52°	3.95°	15.94°	9.02°	5.02°	15.68°	15.76°
	<b>Position</b>	2.63 mm	2.97 mm	5.02 mm	3.78 mm	3.36 mm	3.02 mm	3.58 mm
	<b>Top angle</b>	11.95°	8.56°	10.17°	42.46°	9.34°	20.29°	20.78°
	<b>Active axis</b>	7.24 %	4.3 %	5.6 %	6.3 %	8.10 %	10.97 %	10.67 %
	<b>Active angles</b>	3.40°	2.95°	2.86°	9.6°	3.63°	4.46°	4.31°
Real	<b>Orientation</b>	7.81°	128.83°	19.10°	10.40°	8.93°	8.44°	<b>7.48°</b>
	<b>Position</b>	6.47 mm	69.40 mm	9.42 mm	7.97 mm	7.61 mm	7.30 mm	<b>6.24 mm</b>
	<b>Top angle</b>	15.92°	85.33°	17.54°	35.27°	16.57°	15.81°	<b>13.82°</b>
	<b>Active angles</b>	8.89°	33.04°	9.27°	17.68°	9.16°	8.97°	<b>8.84°</b>

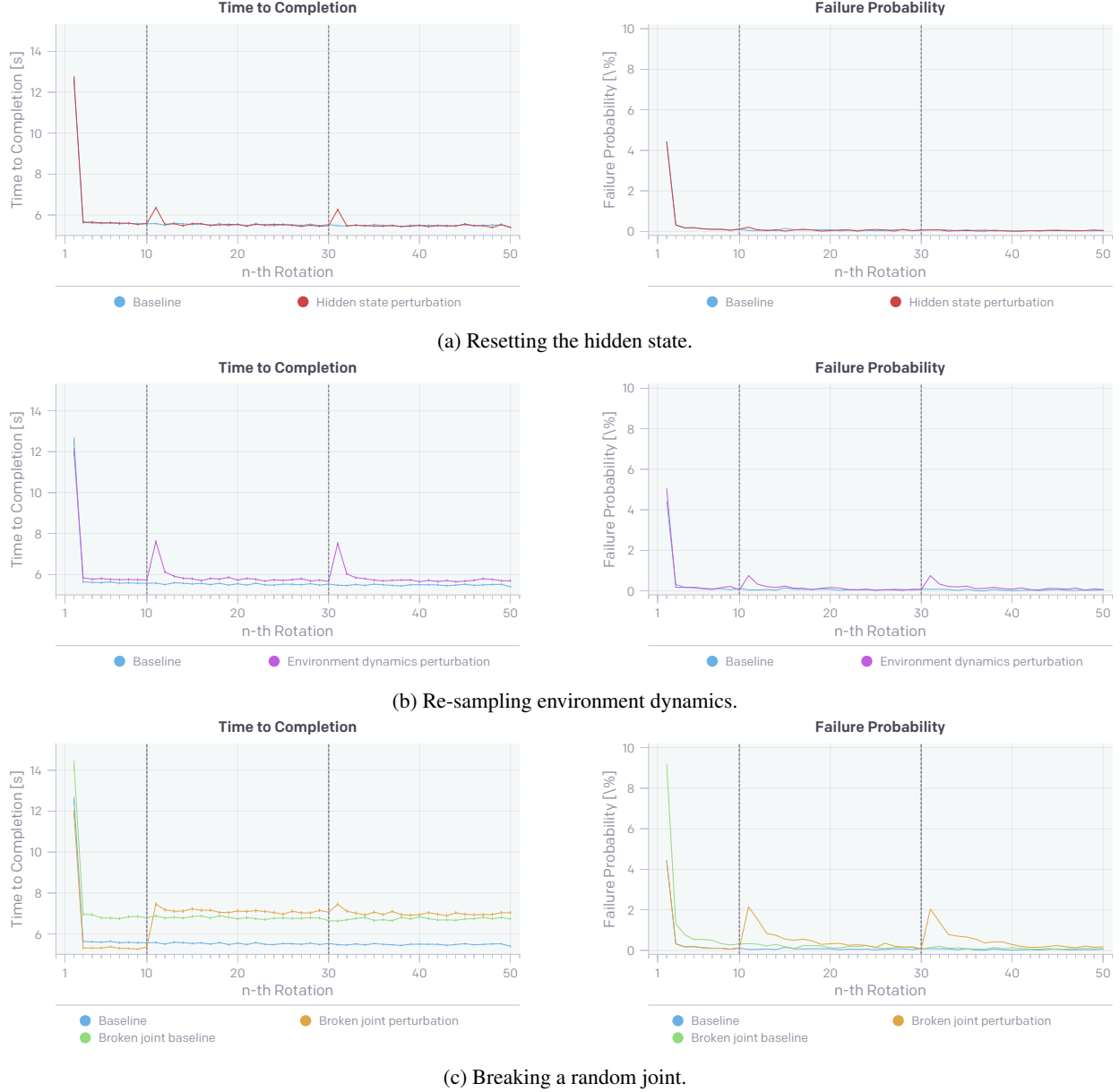


Figure 21: We run 10 000 simulated trials with only face rotations until 50 rotations have been achieved. For each of the face rotations (i.e. the 1<sup>st</sup>, 2<sup>nd</sup>, ..., 50<sup>th</sup>), we measure the average time to completion (in seconds) and average failure probability over those 10k trials. Error bars indicate the estimated standard error. “Baseline” refers to a run without any perturbations applied. “Broken joint baseline” refers to trials where a joint was randomly disabled from the very beginning. We then compare against trials that start without any perturbations but are perturbed at the marked points after the 10<sup>th</sup> and 30<sup>th</sup> rotation by (a) resetting the policy hidden state, (b) re-sampling environment dynamics, or (c) breaking a random joint.



## E Visualizations

### E.1 Policy Recurrent States

We visualize the LSTM cell and hidden states of a Rubik’s cube policy during rollouts. The techniques we use have been successfully applied to study the interpretability of deep neural networks used for hand-writing generation and vision [17, 77].

We collected data from a Rubik’s cube policy rollout in simulation. The data spans 1500 time steps (2 minutes) and contains a sequence of successful cube re-orientations and face rotations. We collected the policy LSTM hidden and cell states at each time step, along with policy observations and actions. Each hidden or cell state vector contained 1024 units.

We arranged the cell (or hidden) state data into a matrix of 1024 rows by 1500 columns. We applied 1D t-SNE [113] to the rows of this matrix in order to group cell state units (“neurons”) that have similar activation patterns across time. We also applied similar tricks as in [17] to further enhance the appearance of the grouped neurons. The resulting visualization is shown in Figure 22 (a).

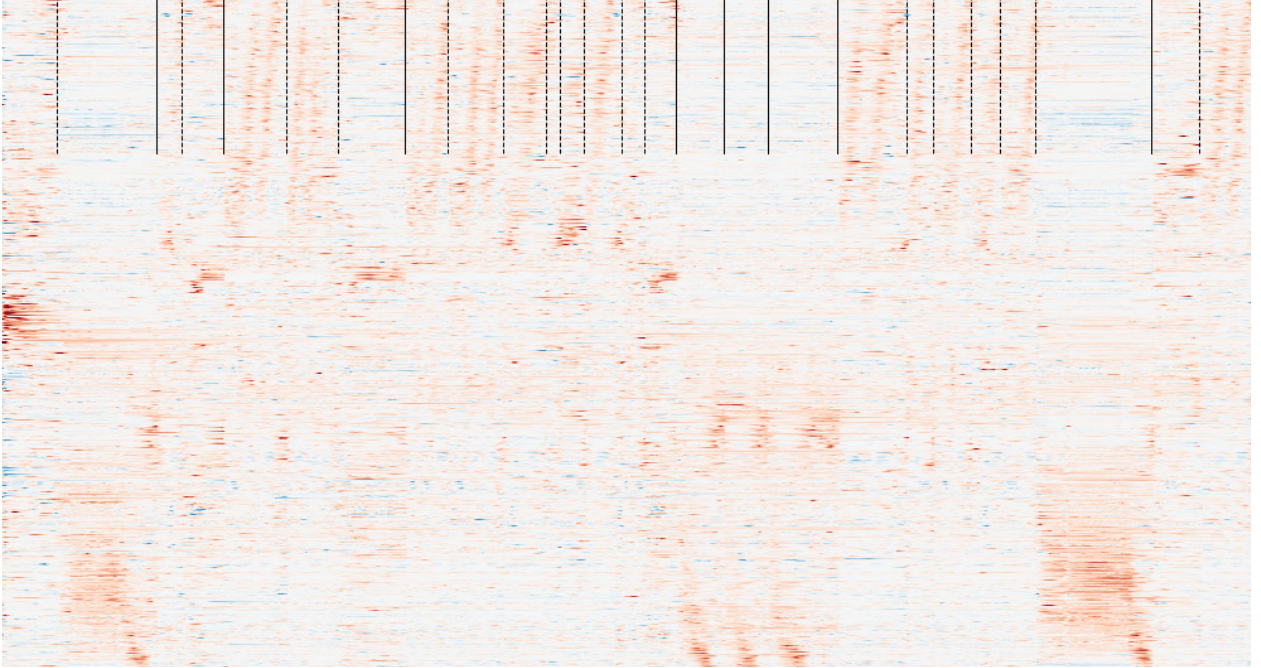
We marked the time steps when the policy successfully completed a block re-orientation (dashed) and a face rotation (solid). We see a clear distinction in activation patterns of groups of neurons in the time before a successful re-orientation vs. face rotation. We can also clearly see when the policy struggles with face rotations (both near the beginning and end of the rollout) where the same group of neurons are active over an extended time.

We performed Non-negative Matrix Factorization (NMF) [77] on the activations matrix to identify factors that can be used to group neurons. We used 6 factors and assigned a distinct color to each factor. We colored each neuron by first projecting the row vector onto each factor then linearly combining the factor colors using the projection weights. The results are shown in Figure 22 (b).

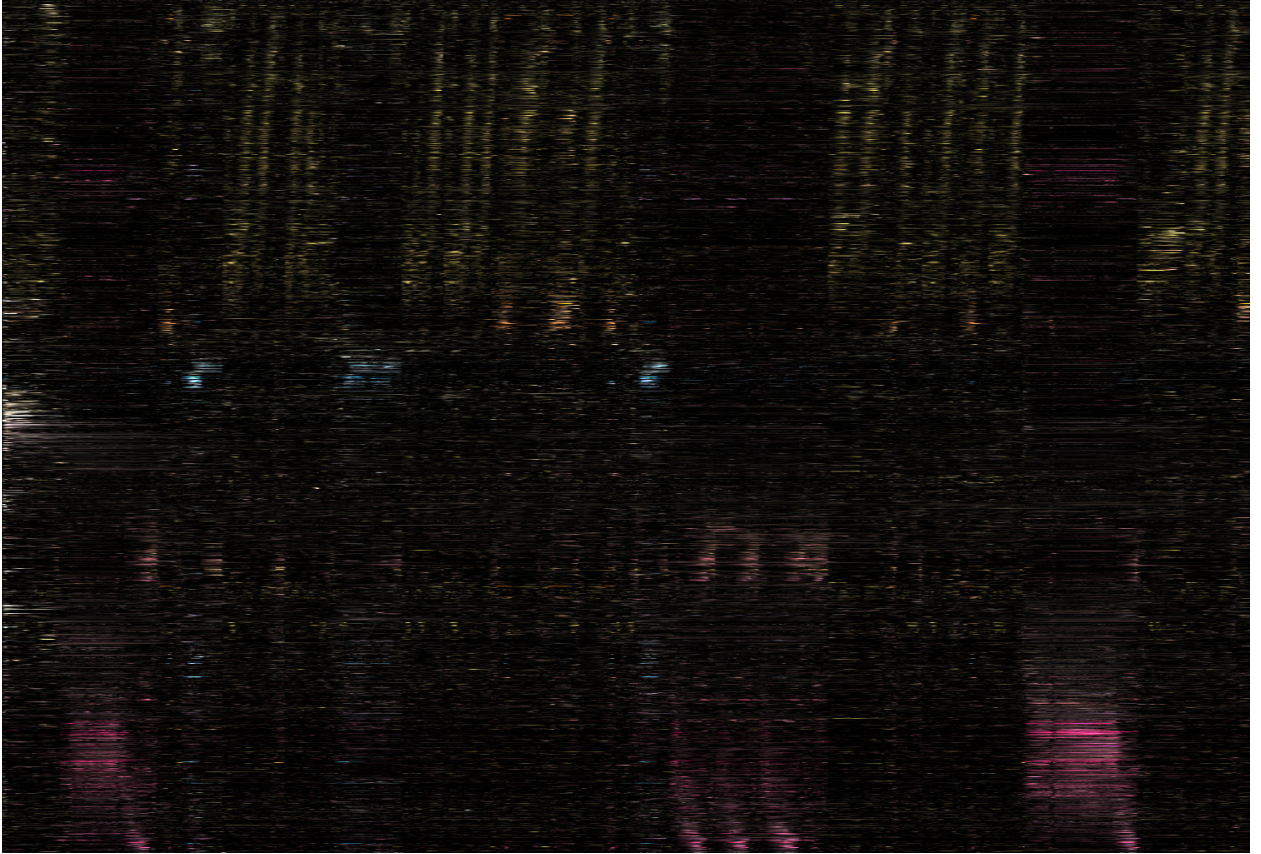
With coloring by NMF factors, we observe 4 distinct groups of neuron activations: yellow, orange, blue, and pink. By matching the activation time of these neuron groups with actions, we found that each group of activations corresponds to a sequence of complex joint actuations, in other words, a skill. For example, the yellow activation group corresponds to flipping the cube back towards the palm of the hand<sup>27</sup>. It is quite surprising to see such high-level skills being clearly identifiable in the cell state representation.

---

<sup>27</sup>See <https://openai.com/blog/solving-rubiks-cube> for the skills corresponding to the other activation groups.



(a)



(b)

Figure 22: Visualization of Rubik's cube policy LSTM cell states. Each row corresponds to the activation values of a cell state (1024 in total) over 1500 time steps. (a) after t-SNE reordering, dashed lines mark a successful block re-orientation, solid lines mark a successful face rotation. (b) after coloring by NMF factors.