

Transferring Dexterous Manipulation from GPU Simulation to a Remote Real-World TriFinger

Arthur Allshire^{1,2}, Mayank Mittal^{2,3}, Varun Lodaya¹, Viktor Makoviychuk², Denys Makoviichuk⁴, Felix Widmaier⁵, Manuel Wüthrich⁵, Stefan Bauer⁶, Ankur Handa², Animesh Garg^{1,2}

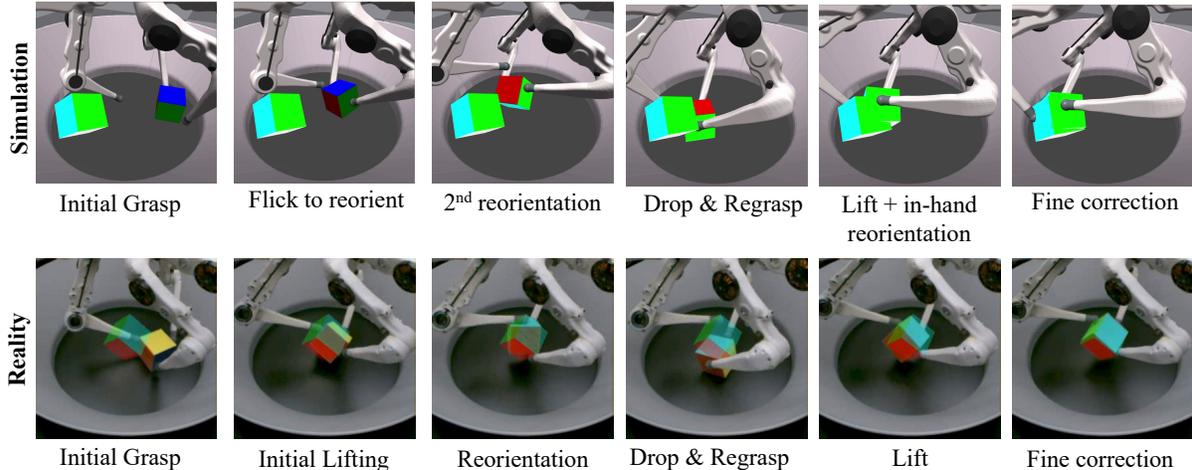


Fig. 1: *Top:* Our system learns to grasp and manipulate objects to 6-DoF goal poses with a single policy, entirely in simulation, across a variety of objects. *Bottom:* We then transfer to a real robot located thousands of kilometers away from where development work is done.

Abstract—In-hand manipulation of objects is an important capability to enable robots to carry-out tasks which demand high levels of dexterity. This work presents a robot systems approach to learning dexterous manipulation tasks involving moving objects to arbitrary 6-DoF poses. We show empirical benefits, both in simulation and sim-to-real transfer, of using keypoint-based representations for object pose in policy observations and reward calculation to train a model-free reinforcement learning agent. By utilizing domain randomization strategies and large-scale training, we achieve a high success rate of 83% on a real TriFinger system, with a single policy able to perform grasping, ungrasping, and finger gaiting in order to achieve arbitrary poses within the workspace. We demonstrate that our policy can generalise to unseen objects, and success rates can be further improved through finetuning. With the aim of assisting further research in learning in-hand manipulation, we provide a detailed exposition of our system and make the codebase of our system available, along with checkpoints trained on billions of steps of experience, at <https://s2r2-ig.github.io>

I. INTRODUCTION

Multi-fingered robotic platforms are essential for executing complicated tasks such as fruit harvesting, and circuit assembly. However, performing such dexterous manipulation tasks requires dealing with various challenging factors. These include high-dimensionality, hybrid dynamics, and uncertainties about the environment [1].

Therefore, performing multi-fingered manipulation in the context of deployed systems poses unique challenges which practical controllers must overcome. In our work, we provide a method for creating controllers capable of achieving reliable

dexterous manipulation tasks across a variety of domain shifts, both from sim to real and across different simulation scenarios. We demonstrate the robustness of our system through: showing transfer of policies created in simulation to the real world, varying system parameters and showing robustness, and changing the object being manipulated.

Our controllers are trained and evaluated on the TriFinger [2] hand. Despite the challenging configuration of the Trifinger system (hand oriented down) and task (grasping and subsequent reposing), we train a unified neural-network to achieve effective and robust control over the system.

The Trifinger robots are run as cloud-based robot farms. Such remote robotic systems promise to alleviate many of the upfront requirements to install and maintain robot hardware [3]. However, RaaS systems are also more rigid than commonly used research platforms, as controllers designed for individual platforms must be rolled out across the entire fleet. Hence, the robustness of our learning-based approach is further proven out by deployment on a system which we lacked physical access to.

The key insight of this paper lies in a careful integration and evaluation of empirical advances in reinforcement learning with high-speed simulation, and practical deployment on a remote robot system. In particular, contribution of this robot systems work are:

- 1) We provide a framework for learning the skill of in-hand manipulation tasks robust to sim-to-real transfer and object morphology.
- 2) Unlike previous in-hand manipulation systems, we produce a single policy which performs both grasping and re-posing, simplifying the pipeline of using RL in such systems.

¹University of Toronto, Vector Institute, ²Nvidia, ³ETH Zurich, ⁴Snap, ⁵MPI Tubingen, ⁶KTH. Email: arthur@allshire.org

- 3) We show the benefits of using keypoints to represent the object in RL algorithms for in-hand manipulation, especially when reposing in 6-DoF.
- 4) We show our system is robust to changing object morphology and physics parameters, shown through simulated experiments and demonstrated sim-to-real transfer.

II. RELATED WORK

A. Learning Dexterous Manipulation with Robot Hands

1. Reinforcement Learning Advances in RL algorithms and computational hardware have enabled rapid progression in the capability of real robots in dynamic scenes. Techniques such as domain randomization and large-scale training have enabled results across a variety of tasks with sim-to-real, including in-hand manipulation [4, 5], as well as in legged locomotion [6, 7]. Active identification of system parameters has also been shown to be helpful in the context of learning manipulation tasks [8].

2. Dexterous Manipulation Dexterous manipulation requires dealing with high-dimensionality of the system, hybrid dynamics, and uncertainties about the environment [1]. Prior work has trained a control policy for in-hand manipulation of a block with a Shadow Dexterous Hand [4]. However, this relied on expensive robot hardware in a controlled environment which could be reset and tuned by engineers. Recent works have extended this setup to other object morphologies [9, 10].

3. Hand Platforms Recently, Wüthrich et al. [2] designed an open-sourced robotic platform for dexterous manipulation called *TriFinger*. The Trifinger platform allows remote deployment of policies on a pre-determined experimental setup [11], without tweaking of system parameters. As a result of the "reset-free" nature of the robot, the hand is facing down, and objects must first be picked up from the ground to be manipulated. Prior works do not provide a method to learn this behaviour in a simple way, either choosing to address only one step in this process or resorting to chaining multiple policies together to accomplish the complete reorientation task [4, 10]. In contrast, in this work, we present an approach to address these shortcomings.

Our system is able to perform 6-DoF in-hand manipulation, as opposed to just 3-DoF reorientation. Furthermore, as the object starts outside of the hand, our single learned policy is able to perform not just picking, but also grasping and un-grasping, in contrast to these prior works.

B. Sim-to-real transfer with RL Policies

1. Challenges of Sim-to-Real Our method relies on learning control policies using gradient-based optimisation combined with large-scale accelerated simulation, a proven method of learning a wide variety of complex robotic tasks [12, 13, 14, 15]. However, doing so means that policies must solve the "sim-to-real transfer" problem of being robust to inference in a different environment than which they were trained. Sim-to-real transfer of manipulation policies is a challenging problem for two major reasons: 1) differences between real-world environment interactions and that of the simulation where policies are trained, and 2) state estimation of objects being manipulated. For the former, a variety of practical methods



(a) OpenAI's shadow hand setup. The cube starts placed in the hand.



(b) The Trifinger setup. The object starts outside of the hand to enable reset-free setup.

Fig. 2: Previous setups for performing RL-based dexterous manipulation in the real world have relied on specialised hardware or configurations which may be impractical to scale. For example, OpenAI's work on Shadow Hand [4] started with the cube in hand (avoiding the need to learn to grasp it), relied on phase space tracking, and only set in-hand orientation goals (rather than full pose goals). In contrast, the Trifinger setup relies only on sensor inputs from RGB cameras and encoders in the fingers, and the object starts in a random position on the ground outside of the hand, yet our system achieves 6-DoF reposing on multiple objects across the workspace.

have been proposed including Domain Randomisation in simulation [15, 14], Bayesian optimisation on the real system [16], and optimisation of the simulator parameters [8, 17].

2. State Estimation Previous methods of state estimation have included pose estimation [4, 18] more recently, policy distillation as a method to solve sim-to-real problem for state estimation [10, 6]. However this is limited by the visual fidelity of current simulators and has only been shown to work for individual objects, limiting the generality of the policy. As a result, practical approaches leveraging existing work in vision still must rely on using pose tracking in the real world.

III. METHOD

A. Problem Setup: The Trifinger Task

1. Task Description. In this paper, we propose a method for training a controller for the Trifinger hand [2] to perform 6-DoF manipulation of objects. The objects start on the ground, and the goal of the system is to move the object to a target position and orientation and hold it there. Any solution must be able to move the fingers to the object, grasp it, and perform appropriate un-grasping and finger-gaiting to achieve the corresponding target orientation.

Our aim is to use Reinforcement Learning (RL) to learn a single policy with a unified reward model to achieve closed-loop controller for this task, in contrast to previous works which have used relied on carefully designed state machines on top of RL based low-level skills to produce such multimodal behaviour (see Fig. 2).

2. Policy Inference on a Remote Real Robot. It is important to show performance on real world rather than simply simulated systems. This is because of the inherent difficulty in having a simulator with the same physics parameters as a real robot, and the necessity of state estimation in the real world, making it the true test of practical methods in robotics.

On the Trifinger system, the pose of the object is tracked on the system using 3 cameras [20]. We convert the position+quaternion representation output by this system into the keypoint representation described in Sec. III-C and use it as input to the policy. Observations of the object pose from the

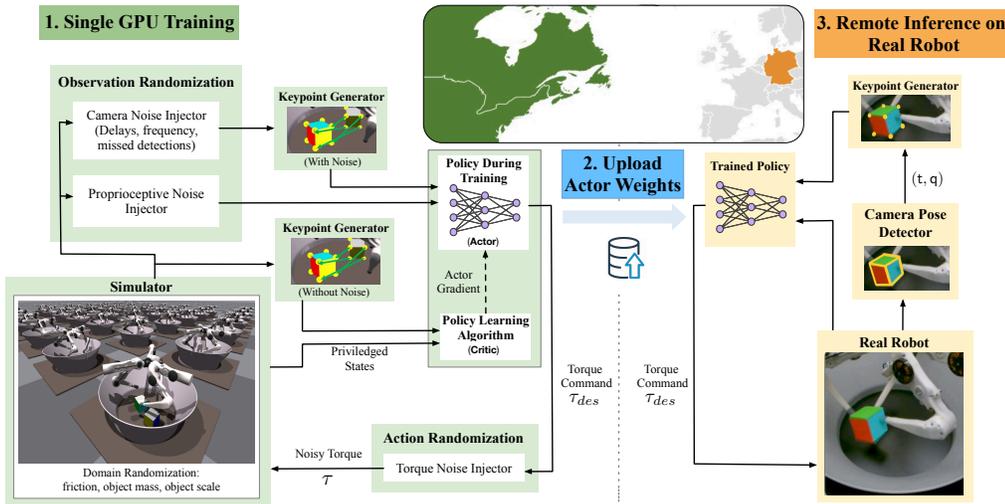


Fig. 3: Our system trains using the IsaacGym simulator¹[12] on 16,384 environments in parallel on a single NVIDIA Tesla V100 or RTX 3090 GPU. Inference is then conducted remotely on a TriFinger robot located across the Atlantic in Germany using the uploaded actor weights. The infrastructure on which we perform sim-to-real transfer is provided courtesy of the organisers of the Real Robot Challenge [19].

OBSERVATION SPACE		DEGREES OF FREEDOM
FINGER JOINTS	POSITION	3 FINGERS · 3 JOINTS · 1 $[\mathbb{R}^1] = 9$
	VELOCITY	3 FINGERS · 3 JOINTS · 1 $[\mathbb{R}^1] = 9$
OBJECT POSE	KEYPOINTS	8 KEYPOINTS · 3 $[\mathbb{R}^3] = 24$
GOAL POSE	KEYPOINTS	8 KEYPOINTS · 3 $[\mathbb{R}^3] = 24$
LAST ACTION	TORQUE	3 FINGERS · 3 JOINTS · 1 $[\mathbb{R}^1] = 9$
TOTAL		75

(a) Actor Observations

OBSERVATION SPACE		DEGREES OF FREEDOM
ACTOR OBSERVATIONS (w/o DR)		75
OBJECT	VELOCITY	6 $[\mathbb{R}^6]$
	POSE	3 FINGERS · 7 $[\mathbb{R}^3 \times SO(3)] = 21$
FINGERTIPS STATE	VELOCITY	3 FINGERS · 6 $[\mathbb{R}^6] = 18$
	WRENCH	3 FINGERS · 6 $[\mathbb{R}^6] = 18$
FINGER JOINTS	TORQUE	3 FINGERS · 3 JOINTS · 1 $[\mathbb{R}^1] = 9$
TOTAL		147

(b) Critic Observations

TABLE I: Asymmetric actor-critic to learn dexterous manipulation. While the actor receives noisy observations, which are added as a part of DR (Sec. III-E), the critic receives the same information without any noise and also has access to certain privileged information from simulator.

camera system are provided at 10Hz, compared to the higher frequency of 50Hz that the policy is run at. This means our method has to deal with relatively low-frequency and noisy object observations based on camera sensing. The torque on each joint is limited such that it does not damage the equipment while in operation, however the exact values of these parameters may vary and are not exposed, meaning our system must be robust to a range of these parameters. The policy is uploaded to the remote system and run on a the robot’s local computer to mitigate latency issues.

B. Reinforcement Learning for Dexterous Manipulation

1. RL Formulation We model our problem using a sequential decision making formulation in which the robotic agent interacts with the environment with the objective of maximising the sum of discounted rewards. This is modelled as a discrete time, partially observable Markov Decision Process (POMDP), represented as the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, P, r, \gamma, \mathcal{S}_0)$, where \mathcal{S} is the state space, \mathcal{O} are the observations corresponding to partial

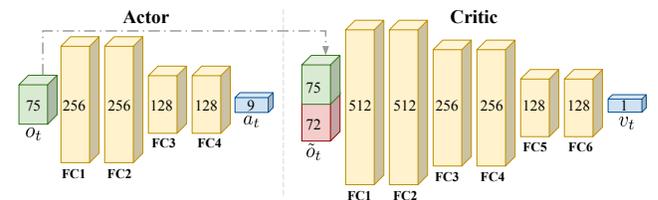


Fig. 4: The actor and critic networks are parameterized using fully-connected layers with ELU activation functions [21].

information about the system states, \mathcal{A} is the action space, $P : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ are the probabilistic state transition dynamics, r is the reward, γ is the discount factor per timestep, and $\mathcal{S}_0 : \mathcal{S} \rightarrow \mathbb{R}$ is the distribution over the system’s initial state at the beginning of an episode.

1. Proximal Policy Optimisation (PPO) is an actor-critic RL algorithm [22] that we build on learning a parametric stochastic policy, $\pi_\theta(a, o)$ mapping from observations to an action distribution to maximise the sum of discounted rewards in each episode. Along with the policy π , PPO learns a value function $V_\phi^\pi(s)$ which approximates the on-policy value function. Following [23], the learned value function is a function of states \mathcal{S} rather than observations \mathcal{O} , which improves the accuracy of value function estimates. The observations $o \in \mathcal{O}$ of the policy are described in Table Ia and the states $s \in \mathcal{S}$ provided in the value function are described in Table Ib.

2. Parametrization While multiple formulations of reward and observations are possible, we choose to use a parametrization based on a keypoint formulation to represent the object pose, which we find boosts the ability of the RL algorithm to learn the task at hand (see Sec. III-C, III-D, IV-B).

3. Hyper-parameters For PPO, we use the following hyper-parameters: discount factor $\gamma = 0.99$, clipping $\epsilon = 0.2$. The learning rate is annealed linearly over the course of training from $5e-3$ to $1e-6$. Our policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ is a Multilayer perceptron (MLP) with 4 hidden layers, 2 of size 256 followed by 2 of size 128, and 9 outputs which are scaled to the torque ranges of the real robot. Our value function $V_\phi^\pi : \mathcal{S} \rightarrow \mathbb{R}$ is an MLP with 2 layers of size 512, followed by 2 layers of size 256 and 128 each and produces a scalar value function

PARAMETER	RANGE	σ	σ_{corr}
OBSERVATION NOISE			
OBJECT POSITION ²	[-0.30, 0.30]	0.002	0.000
OBJECT ORIENTATION ²	[-1.00, 1.00]	0.020	0.000
FINGER JOINT POSITION	[-2.70, 1.57]	0.003	0.004
FINGER JOINT VELOCITY	[-10.00, 10.0]	0.003	0.004
ACTION NOISE			
APPLIED JOINT TORQUE	[-0.36, 0.36]	0.02	0.01
PARAMETER	SCALING DISTRIBUTION		
ENVIRONMENT PARAMETERS			
OBJECT SCALE	UNIFORM(0.97, 1.03)		
OBJECT MASS	UNIFORM(0.70, 1.30)		
OBJECT FRICTION	UNIFORM(0.70, 1.30)		
TABLE FRICTION	UNIFORM(0.50, 1.50)		
EXTERNAL FORCES	REFER TO [4, PP. 9]		

TABLE II: For observations and actions, σ and σ_{corr} are the standard deviation of additive gaussian noise sampled every timestep and at the start of each episode, respectively. For environment, the parameters represent scaling factor applied to the nominal values in the real robot model.

as output. The action space \mathcal{A} of our policy is torque on each of the 9 joints of the robot.

4. RL Library We build on the implementation of PPO from the RL Games library [24], which vectorizes observations and actions on GPU allowing us to take advantage of the parallelising provided by the simulator (see Sec. III-E) by reducing the overhead in CPU-GPU communication typical to most CPU-based simulators.

C. Representation of the Object Pose: Keypoints

We focus on the task of manipulating an object in 6 degrees of freedom. As such, we must represent the pose of the object at multiple stages of our training pipeline. In order to capture both position and orientation in the same space in our representation, we use eight keypoints at the edges of the oriented bounding box of the object being manipulated. In the object’s local frame these are denoted $k_i^l \in \mathbb{R}^3$, $i = 1, \dots, 8$. The keypoints in the world frame are related to those in the local frame by a transformation $k_i^c = \mathbf{T}k_i^l$, where \mathbf{T} depends on the current pose of the object.

In Sec. IV-B, we contrast keypoint representations to a position+quaternion formulation used in [4, 25], finding that keypoints improve the policy’s success rate substantially. During policy inference in the real world, we note as long as we are able to get the bounding box of the object (via classical trackers as on the Trifinger [20], or with learning-based setups which commonly rely on the same keypoints on the oriented bounding box, such as [18, 26]), we are able to obtain the keypoints on the object bounding box required to use this representation in policy input.

D. Reward Formulation & Curriculum

1. Kernel Our reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow R$ has three components. Following [6], we use a logistic kernel to convert tracking error in Euclidean space into a bounded reward function. We generalise the kernel formulation to account for a range of distance scales, defining, $\mathcal{K}(x) =$

²The noise to keypoints is not applied directly, instead it is added to the object pose in the world frame before computing the keypoints through it.

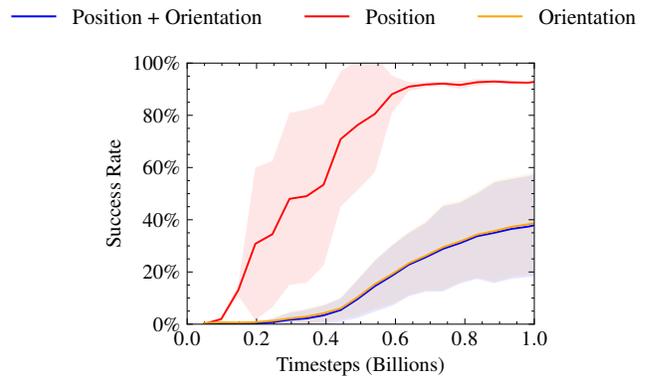


Fig. 5: Training curves on a reward function similar to prior work [16, 28] for the setting with DR. We take the average of 5 seeds; the shaded areas show standard deviation, noting that curves for Orientation and Position+Orientation overlap during training. It is worth noting that the nature of the reward makes it very difficult for the policy to optimize, particularly achieving an orientation goal.

$(e^{ax} + b + e^{-ax})^{-1}$, where a is a scaling factor and b controls the sensitivity of the kernel at low distances.

2. Object Displacement Reward As noted in Sec. III-C, we use keypoints in order to calculate the reward in a natural space for 3-D reposing. The component of the reward corresponding to the difference between the object’s current pose and the desired target pose is given by $r_o = \sum_{i=1}^N \mathcal{K}(\|k_i^c - k_i^T\|)$, where the k_i^c and k_i^T lie at the $N = 8$ vertices of the bounding boxes of the object at the current and target configurations respectively (see Sec. III-C).

3. Finger Reaching Reward To encourage the fingers to reach the object during initial exploration, we give a reward for moving the fingers towards the object, which was also found to be helpful in [27]. This term is defined by sum of the movement of each fingertip towards the goal per timestep: $r_f = \sum_{i=1}^3 \Delta_i^t$, where Δ denotes the change across the timestep of the fingertip distance to the centroid of the object, $\Delta_i^t = \|\mathbf{f}_{i,t} - \mathbf{p}_t^c\|_2 - \|\mathbf{f}_{i,t-1} - \mathbf{p}_{t-1}^c\|_2$, and $\mathbf{f}_i \in \mathbb{R}^3$ denotes the position of the i -th fingertip, and \mathbf{p}_t^c denotes the position of the centroid of the object.

Finally, we define a penalty on the movement of each finger: $r_v = \sum_{i=1}^3 \|\dot{\mathbf{f}}_i\|_2^2$ where $\dot{\mathbf{f}}_i$ denotes the velocity of the i ’th fingertip in the global frame.

4. Total Reward Our total reward is defined as:

$$R(s, a) = w_f \times r_f \times \mathbf{I}(t \leq N_v) + w_v \times r_v + w_o \times r_o$$

where $w_f = -750$, $w_v = -0.5$ and $w_o = 40$ are the weights of each reward component, determined through search over numerous training runs. We also found in initial experimentation that the curriculum reducing the weight of r_f reward to 0 after $N_v = 5e7$ timesteps was needed in order to allow the robot to perform ungrasping needed to facilitate reorientation which is learned later in training. However, having the reward term during the initial phases of training sped up learning by encouraging the robot to interact with the object. In the kernel \mathcal{K} we use $a = 30$ and $b = 2$ which provided a good balance between learning behaviour early in training and good accuracy later in training.

E. Simulation Environment

1. Choice of Simulator We train on the IsaacGym simulator [12], a simulation environment tailored towards allowing policy learning with a high sampling rate by paralling physics on a single GPU (>50K samples/sec in policy inference on Tesla V100 and around 100K samples/sec on RTX 3090). A high sampling rate is essential to learn complex dynamic robotics tasks quickly [12, 13], and the ability to perform training and inference on desktop-level systems is important for enabling other researchers to build on our work and use it for in-hand manipulation.

2. Domain Randomisation Domain Randomization (DR) is a method for improving the robustness of policies for sim-to-real transfer [29, 15, 14]. In Domain Randomisation, simulator parameters are sampled from a distribution, $\xi \sim P(\Xi)$ in order to modify the physics behaviour of the simulated environment. If the real-world physics parameters ξ are within the support of the distribution $P(\Xi)$, then a policy which successfully achieves the task in simulation will be able to perform comparably in the real world.

We choose our Domain Randomization parameters to account for modelling errors in the environment as well as noise in sensor measurements. These parameters are listed in Table II. In addition to these randomizations, we apply random forces to the object as described in [4] in order to improve the stability of grasps and represent un-modelled dynamics. We mimic the refresh rate of the camera on the real system, by repeating the observation of the keypoints for 5 frames (See Sec. III-E). To mimic possible extra camera latency, with 3% probability, we repeat the camera-based object-pose observations for subsequent rounds of policy to mimic dropped frames from the tracker. Up-to-date proprioceptive data is provided to allow the policy to take advantage of the high-frequency and more reliable encoder information available on the real system.

IV. EXPERIMENTS

In our experiments, we aim to answer the following four questions pertaining to learning a robust policy for this task, as well as evaluating how well it transfers to the real world:

- 1) How well does our system learn 6-DoF manipulation with a reward function based on prior works?
- 2) How does performance change when we use a task-appropriate representation - keypoints - for reward computation and policy input in the 6-DoF reposing task?
- 3) Is our system robust to sensor noise and varying environment parameters, and to changes in object morphology?
- 4) How well do our policies, trained entirely in simulation, transfer to the real TriFinger system?

A. Experiment 1: Training

1. Success Criterion The aim in our 6-DoF manipulation task is to get the position and orientation of the object to a specified goal position and orientation. We define our metric for ‘success’ in this task as getting the position within 2 cm, and orientation within 0.4 rad (22°) of the target goal pose as used in [4]; comparable to mean results obtained in [16]. Following previous works dealing with similar tasks [4, 5, 27],

we apply a reward based on the position and orientation components of error individually.

2. Alternative Reward Formulation Following experimentation, the best candidate reward of this format was:

$$r_o = \mathcal{K}(\|\mathbf{t}^C - \mathbf{t}^T\|_2) + \frac{1}{3 \times |d^r| + 0.01}$$

where \mathbf{t}^C and \mathbf{t}^T are the current and goal positions of the object, $d^r = 2 \times \arcsin(\min(1.0, \|\mathbf{q}_{\text{diff}}\|_2))$, $\mathbf{q}_{\text{diff}} = \mathbf{q}^C(\mathbf{q}^T)^*$. \mathcal{K} is the logistic kernel that takes L2 norm between the current and target object position as input, and d^r is the distance in radians between the current and target object orientation. We use the alternative scaling parameter $a = 50$ in \mathcal{K} , which we found to work better in this reward formulation (see Sec. III-D). We use the same weightings for each of the 3 components of the reward as in Sec. III-D. For this experiment, we trained on the 6.5cm cube used in the Real Robot Challenge [19].

3. Results The results are shown in Fig. 5. We found that while this formulation of the reward was good at allowing PPO to learn a policy to get the object to the goal, even after 1 Billion steps in an environment with no Domain Randomization it was learning very slowly to achieve the orientation goal.

B. Experiment 2: Representation of Pose

1. Comparison The poor results in Experiment 1 (Sec. IV-A) lead us to search for alternative representations of object pose in the calculation of the reward and policy observations; these are described in Sec. III-C and Sec. III-D. We compared our method of using keypoints to represent the object pose and using positions and quaternions in two ways: firstly, using it as the policy input as compared to a position and quaternion representation, and secondly, using it to calculate the reward as compared to a reward based on the linear and angular rotational distances individually.

2. Experimental Setup For the observations, in order to provide a fair comparison between position/quaternion and keypoints as policy input, observation noise and delays are applied in the same manner (by applying them in the position and quaternion space before transforming to keypoints, as noted in Sec. III-E). Also note that both representations only rely on the spatial pose information and fixed size of the object to compute. The pose of the object is represented with a 7-dim vector involving translation and quaternion (\mathbf{t} , \mathbf{q}). The position and quaternion of the goal pose are provided as input to the actor and critic, replacing the keypoints in Tables Ia and Ib.

For the reward, in order to provide a fair comparison to the keypoints reward, as mentioned previously, many hours were spent tuning the kernels and parameters used in the translation based reward, described in Experiment 1. In comparison, little effort was spent tuning the keypoints function, with only one tweak to the weightings in the logistic kernel, showing the relative simplicity of working with this formulation. For this experiment, we trained on the 6.5cm cube used in the Real Robot Challenge [19], with keypoints placed at the bounding box (in this case the corners of the cube).

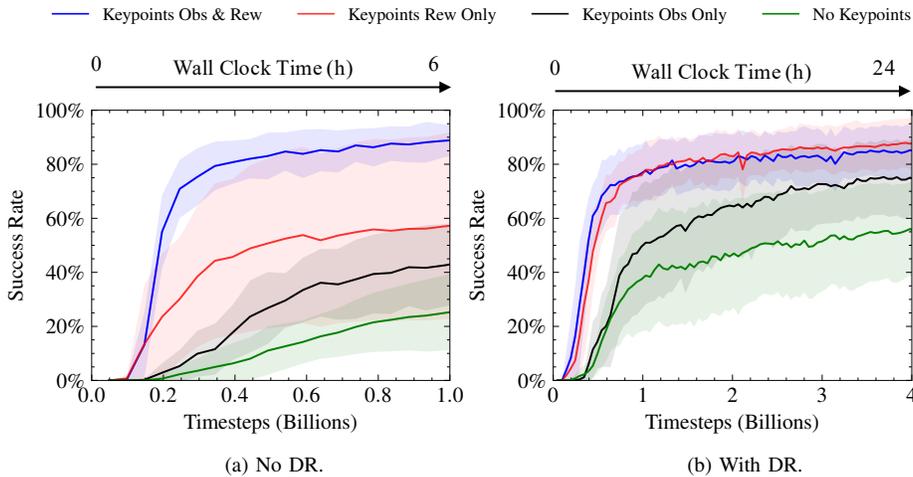


Fig. 6: Success Rate over the course of training without and with domain randomization. Each curve is the average of 5 seeds; the shaded areas show standard deviation. Note that training without DR is shown to 1B steps to verify performance; use of DR didn't have a large impact on simulation success rates after initial training.

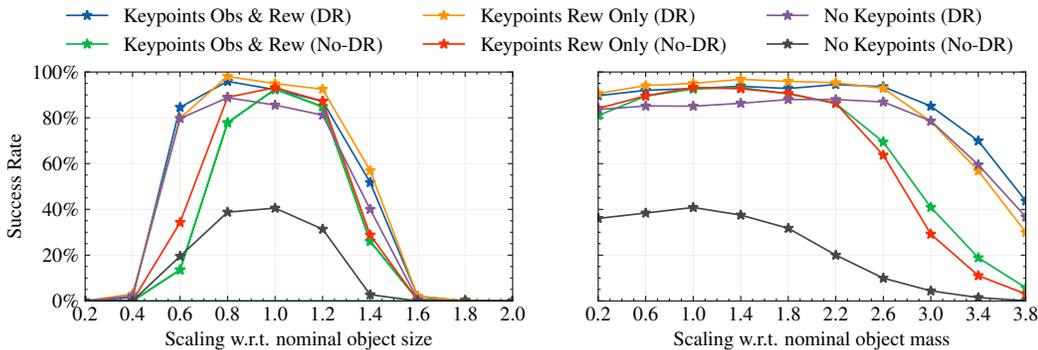


Fig. 7: We show the robustness to varying the object parameters outside the DR range it was trained for. In the evaluations, all the other DR is turned off to ensure a controlled setting. Each success is evaluated over 1024 runs with random goal and object initialization.

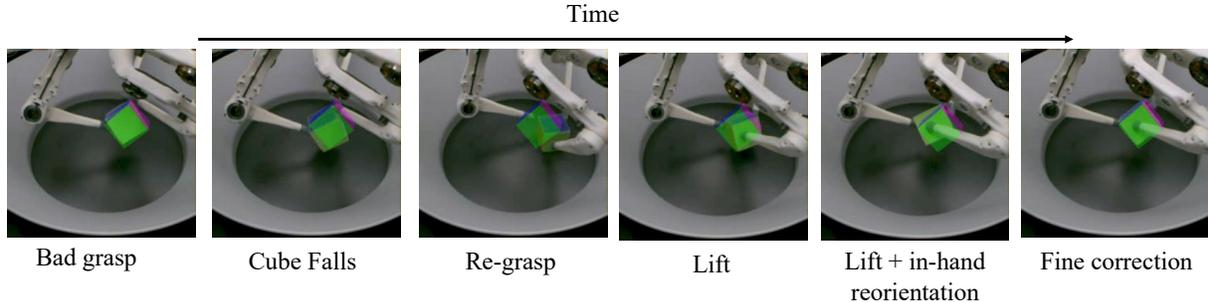


Fig. 8: The use of a single, continuous policy for grasping and reposing allows the policy to automatically recover from failures. For example in the sequence above the system recovers from a failure and re-grasps the cube to achieve the desired goal pose.

3. Results Fig. 6 shows the results of training, with both timesteps and wall-clock time. In the curve without any Domain Randomization, we trained for 1 billion steps over the course of 6 hours on a single GPU. Using keypoints in observations and the reward function performs the best of the four policies, also exhibiting a low variance among seeds.

When Domain Randomization is applied, the two curves with a keypoints-based reward are far better in terms of success rate at the end of training and in terms of convergence rate; however, in this case having keypoint observations seems to matter somewhat less. This is perhaps due to the longer training (4b steps & 24 hours on a single GPU) overwhelming the inductive bias introduced by using keypoints as representations. However, using keypoints to compute the reward provided a large benefit in both cases,

showing the improvement caused by calculating the reward in Euclidean space rather than mixing linear and angular displacements through addition.

C. Experiment 3: Robustness of Policies in Simulation

1. Impact of Varying Physics Parameters In order to investigate the impact that Domain Randomization (see Sec. III-E) has on the robustness of policies of a hand in this configuration, we ran experiments by varying parameters outside of the normal domain randomization ranges in simulation. Fig. 7 shows the results. We find that, despite only being randomized initially within a range of 0.97-1.03x nominal size, our policies with Domain Randomization achieve over an 80% success rate even with a scale of 0.6 and 1.2x nominal size, while those without DR have a success rate that drops off much more quickly outside the normal range. We find similar results when scaling the object mass

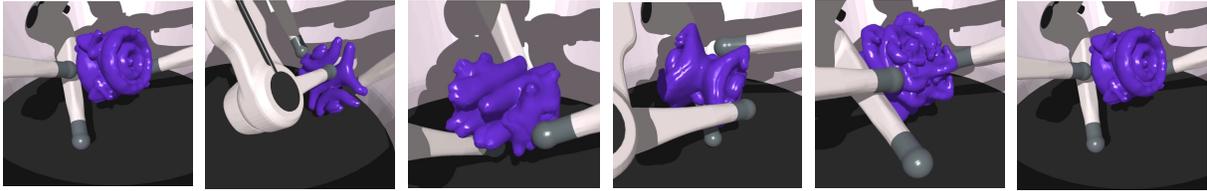


Fig. 9: A selection of the trifinger manipulating various different objects from the EGAD dataset. The system achieves a diversity of emergent grasps enabling it to manipulate a variety of object morphologies to 6-DoF target poses. See the website (<https://s2r2-ig.github.io>) for more videos of object manipulation.

Object	Cube Policy	Finetuned
Cube 6.5cm ³	92.1 %	-
EGAD (train)	84.9 %	87.2 %
EGAD (test)	85.2 %	86.4 %
Different Sized Cuboids	46.8 %	72.9 %

TABLE III: Fine-tuning performance. We tested the zero-shot performance on the EGAD dataset and a number of differently sized cuboids to measure the robustness to differing object morphologies and scales, respectively. Numbers calculated from N=1024 trials in simulation. We disabled environment, observation & action randomizations.

relative to the nominal range, however in this case we find that the policies using keypoints-based reward even without DR is much more robust at masses 3x nominal.

2. Other objects & Fine-Tuning We tested the change in performance resulting from finetuning the policy (O-KP+R-KP trained with DR on the 6.5cm³ cube) on 100 randomly-selected objects from the EGAD dataset, and on a selection of cuboids with random side lengths between 2-8 cm. We found an improvement in policy performance resulting from finetuning on 100 randomly selected EGAD objects, on both this same set of objects and test objects (see Table III). The gains were even bigger when fine-tuning on cuboids of different scales, suggesting that scale is a more important consideration than shape for generalisation.

D. Experiment 4: Simulation to Remote Real Robot Transfer

1. Experimental Setup We ran experiments on the real robot to determine the success rate of the policies trained with Domain Randomization under the metric defined in Sec. 5. We performed $N = 40$ trials for each policy on the task setup described in Sec. III-E; the results for each of the four ablations on keypoints which we tested are shown in Fig. 10.

2. Results Out of the four models discussed in Sec. IV-B, the best policy achieved a success rate of 82.5%. This was achieved with the use of keypoints used in observations of the policy as well as the reward function during training (O-KP+R-KP). The policy using position+quaternion representations but with a reward calculated with keypoints (O-PQ+R-KP) achieved a 77.5% success rate. These first two policies were well within each others' confidence intervals. This is likely due to the impact of the better representation of keypoints being mitigated somewhat after 4 Billion steps of training, as discussed in Sec. IV-B. In contrast, neither of the policies trained using the position & quaternion based reward achieved good success rates, with the policy using keypoints-based observations (O-KP+R-PQ) achieving only a 60% success rate while the one with position and quaternion observations (O-PQ+R-PQ) only achieved a 55% success rate. These results

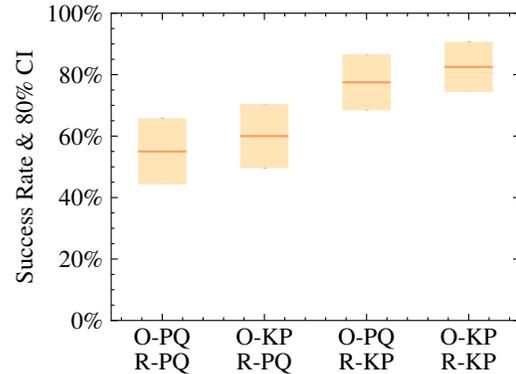


Fig. 10: Success Rate on the real robot plotted for different trained agents. O-PQ and O-KP stand for position+quaternion and keypoints observations respectively, and R-PQ and R-KP stand for linear+angular and keypoints based rewards respectively, as discussed in Sec. IV-B. Each mean made of N=40 trials and error bars calculated based on an 80% confidence interval.

show the importance of having a reward function which effectively balances learning to achieve the goal in \mathbb{R}^3 and $SO(3)$ in order to have policies with a high success rate in simulation, and thus a high corresponding success rate after real robot transfer.

3. Qualitative Behaviour We noticed a variety of emergent behaviours used to achieve sub-goals within the overall object-reposing task. We display some of these in the panel in Figures 1 and 8. The most prominent of these is "dropping and regrasping". In this maneuver, the robot learns to drop the cube when it is close to the correct position, re-grasp, and pick it back up. This enables the robot to get a stable grasp on the cube in the right position. The robot learns to use the motion of the object to the correct location in the arena as an opportunity to simultaneously rotate it on the ground to make achieving the correct grasp in challenging target locations far from the center of the fingers' workspace. Our policy is also robust to dropping - it can recover from an object falling out of the hand and retrieve it from the ground.

We were only able to perform these experiments with the 6.5cm³ cube, as the Trifinger remote inference setup only provided a single size of cube to test on. However, using an off the shelf pose detector (for example [18, 26]), sim-to-real with these is a direction that we could pursue in the future.

V. SUMMARY

This paper emphasizes the empirical value of a systems approach to robot learning through a case study in dexterous manipulation. We introduced a framework for learning in-hand manipulation tasks and transferring the resulting policies to the real world. We show how RL algorithms for in-hand

manipulation can benefit from using keypoints as opposed to the more ordinary angular and linear displacement-based reward and observation computation. We show that our policies are able to generalise to unseen objects, and success rates can be further improved through finetuning. In contrast to prior work, our system solves all of the challenges inherent in 6-DoF grasping and reposing in a single policy, simplifying the pipeline of using RL for dexterous manipulation. We provide a clear elucidation of our approach and open source checkpoints and code to allow reproducing our work.

REFERENCES

- [1] A. M. Okamura, N. Smaby, and M. R. Cutkosky, "An overview of dexterous manipulation," in *Proceedings 2000 ICRA*. IEEE, 2000.
- [2] M. Wüthrich, F. Widmaier, F. Grimmering, J. Akpo, S. Joshi, V. Agrawal, B. Hammoud, M. Khadiv, M. Bogdanovic, V. Berenz, J. Viereck, M. Naveau, L. Righetti, B. Schölkopf, and S. Bauer, "Trifinger: An open-source robot for learning dexterity," 2020. [Online]. Available: <https://arxiv.org/abs/2008.03596>
- [3] B. Kehoe, S. Patil, P. Abbeel, and K. Goldberg, "A survey of research on cloud robotics and automation," *IEEE Transactions on Automation Science and Engineering*, no. 2, 2015.
- [4] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. W. Pachocki, 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," 2018. [Online]. Available: <http://arxiv.org/abs/1808.00177>
- [5] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, "Solving rubik's cube with a robot hand," 2019. [Online]. Available: <http://arxiv.org/abs/1910.07113>
- [6] 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*, no. 26, Jan 2019.
- [7] F. Shi, T. Homberger, J. Lee, T. Miki, M. Zhao, F. Farshidian, K. Okada, M. Inaba, and M. Hutter, "Circus anymal: A quadruped learning dexterous manipulation with its limbs," 2020.
- [8] 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," 2019.
- [9] W. Huang, I. Mordatch, P. Abbeel, and D. Pathak, "Generalization in dexterous manipulation via geometry-aware multi-task learning," 2021. [Online]. Available: <https://arxiv.org/abs/2111.03062>
- [10] T. Chen, J. Xu, and P. Agrawal, "A system for general in-hand object re-orientation," 2021. [Online]. Available: <https://arxiv.org/abs/2111.03043>
- [11] S. Bauer, F. Widmaier, M. Wüthrich, N. Funk, J. U. D. Jesus, J. Peters, J. Watson, C. Chen, K. Srinivasan, J. Zhang, J. Zhang, M. R. Walter, R. Madan, C. B. Schaff, T. Maeda, T. Yoneda, D. Yarats, A. Allshire, E. K. Gordon, T. Bhattacharjee, S. S. Srinivasa, A. Garg, A. Buchholz, S. Stark, T. Steinbrenner, J. Akpo, S. Joshi, V. Agrawal, and V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and G. State, "Isaac Gym: High Performance GPU-Based Physics Simulation For Robot Learning," 2021. [Online]. Available: <https://arxiv.org/abs/2108.10470>
- [12] B. Schölkopf, "A robot cluster for reproducible research in dexterous manipulation," vol. abs/2109.10957, 2021. [Online]. Available: <https://arxiv.org/abs/2109.10957>
- [13] C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem, "Brax - a differentiable physics engine for large scale rigid body simulation," 2021. [Online]. Available: <http://github.com/google/brax>
- [14] A. Mandlekar*, Y. Zhu*, A. Garg*, L. Fei-Fei, and S. Savarese (* equal contribution), "Adversarially Robust Policy Learning through Active Construction of Physically-Plausible Perturbations," in *IROS 2017*, 2017.
- [15] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *ICRA 2018*, May 2018.
- [16] N. Funk, C. B. Schaff, R. Madan, T. Yoneda, J. U. D. Jesus, J. Watson, E. K. Gordon, F. Widmaier, S. Bauer, S. S. Srinivasa, T. Bhattacharjee, M. R. Walter, and J. Peters, "Benchmarking structured policies and policy optimization for real-world dexterous object manipulation," 2021. [Online]. Available: <https://arxiv.org/abs/2105.02087>
- [17] F. Ramos, R. C. Possas, and D. Fox, "Bayessim: adaptive domain randomization via probabilistic inference for robotics simulators," 2019. [Online]. Available: <http://arxiv.org/abs/1906.01728>
- [18] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, "Deep Object Pose Estimation for Semantic Robotic Grasping of Household Objects," in *CoRL*, 2018.
- [19] "Real robot challenge," 2020. [Online]. Available: <https://real-robot-challenge.com>
- [20] M. Wüthrich, F. Widmaier, O. Ahmed, and V. Agrawal, "Trifinger object tracking," 2020. [Online]. Available: https://github.com/open-dynamic-robot-initiative/trifinger_object_tracking
- [21] D. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," in *ICLR 2016*, Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: <http://arxiv.org/abs/1511.07289>
- [22] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017.
- [23] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning," 2017. [Online]. Available: <http://arxiv.org/abs/1710.06542>
- [24] D. Makoviychuk and V. Makoviychuk, "Rl games," 2021. [Online]. Available: https://github.com/Denys88/rl_games/
- [25] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox, "GPU-Accelerated Robotic Simulation for Distributed Reinforcement Learning," in *CoRL 2018*. PMLR, 2018. [Online]. Available: <http://proceedings.mlr.press/v87/liang18a.html>
- [26] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic, "Cosypose: Consistent multi-view multi-object 6d pose estimation," 2020. [Online]. Available: <https://arxiv.org/abs/2008.08465>
- [27] O. Ahmed, F. Träuble, A. Goyal, A. Neitz, M. Wüthrich, Y. Bengio, B. Schölkopf, and S. Bauer, "Causalworld: A robotic manipulation benchmark for causal structure and transfer learning," *CoRR*, 2020. [Online]. Available: <https://arxiv.org/abs/2010.04296>
- [28] C. Chen, K. Srinivasan, J. Zhang, and J. Zhang, "Dexterous manipulation primitives for the real robot challenge," *CoRR*, 2021. [Online]. Available: <https://arxiv.org/abs/2101.11597>
- [29] 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," *CoRR*, 2017. [Online]. Available: <http://arxiv.org/abs/1703.06907>