
A Framework for Efficient Robotic Manipulation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent advances in unsupervised representation learning significantly improved
2 the sample efficiency of training Reinforcement Learning policies in simulated
3 environments. However, similar gains have not yet been seen for real-robot learning.
4 In this work, we focus on enabling data-efficient real-robot learning from pixels. We
5 present a **Framework for Efficient Robotic Manipulation (FERM)**, a method that
6 utilizes data augmentation and unsupervised learning to achieve sample-efficient
7 training of real-robot arm policies from sparse rewards. While contrastive pre-
8 training, data augmentation, and demonstrations are alone insufficient for efficient
9 learning, our main contribution is showing that the combination of these disparate
10 techniques results in a simple yet data-efficient method. We show that, given
11 only 10 demonstrations, a single robotic arm can learn sparse-reward manipulation
12 policies from pixels, such as reaching, picking, moving, pulling a large object,
13 flipping a switch, and opening a drawer in just 30 minutes of mean real-world
14 training time.

15 1 Introduction

16 Recent advances in deep reinforcement learning (RL) have given rise to unprecedented capabilities in
17 autonomous decision making. Notable successes include learning to solve a diverse set of challenging
18 video games Mnih et al. (2015); Berner et al. (2019); Vinyals et al. (2019); Badia et al. (2020),
19 mastering complex classical games like Go, Chess, Shogi, and Hanabi Silver et al. (2016, 2017);
20 Schrittwieser et al. (2019), and learning autonomous robotic control policies in both simulated
21 Schulman et al. (2015, 2017); Laskin et al. (2020); Hafner et al. (2020) and real-world settings Levine
22 et al. (2015); Kalashnikov et al. (2018). In particular, deep RL has been an effective method for
23 learning diverse robotic manipulation policies such as grasping Pinto and Gupta (2016); Mahler
24 et al. (2016); Levine et al. (2016); Gupta et al. (2018) and dexterous in-hand manipulation of
25 objects Andrychowicz et al. (2018).

26 However, to date, general purpose RL algorithms have been extremely sample inefficient, which has
27 limited their widespread adoption in the field of robotics. State-of-the-art RL algorithms for discrete
28 Hessel et al. (2017) and continuous Lillicrap et al. (2015) control often require tens of millions of
29 environment interactions to learn effective policies from image input Tassa et al. (2018), while training
30 the Dota5 agent Berner et al. (2019) to perform competitively to human experts required an estimated
31 180 human-years of game play. Even when the underlying proprioceptive state is accessible, sparse
32 reward robotic manipulation still needs millions of training samples Andrychowicz et al. (2017), an
33 estimated 2 weeks of training in real time, to achieve reliable success rates on fundamental tasks such
34 as reaching, picking, pushing, and placing objects.

35 Another common approach to learned robotic control is through imitation learning Zhang et al. (2018);
36 Ho and Ermon (2016); Duan et al. (2017); Finn et al. (2017); Young et al. (2020), where a large
37 number of expert demonstrations are collected and the policy is extracted through supervised learning
38 by regressing onto the expert trajectories. However, imitation learning usually needs hundreds or

39 thousands of expert demonstrations, which are laborious to collect, and the resulting policies are
40 bounded by the quality of expert demonstrations. It would be more desirable to learn the optimal
41 policy required to solve a particular task autonomously.

42 In this work, rather than relying on transferring policies from simulation or la-
43 bor intensive human input through imitation learning or environment engineering,
44 we investigate how pixel-based RL applied to
45 real robots can be made data-efficient. Recent
46 progress in unsupervised representation learn-
47 ing Laskin et al. (2020); Stooke et al. (2020)
48 and data augmentation Laskin et al. (2020);
49 Kostrikov et al. (2020) has significantly im-
50 proved the efficiency of learning with RL in
51 simulated robotic Tassa et al. (2018) and video
52 game Bellemare et al. (2013) environments. The
53 primary strength of these methods is learning
54 high quality representations from image input
55 either explicitly through unsupervised learning
56 or implicitly via data augmentation.

57 Building on these advances, we propose
58 **Framework for Efficient Robotic Manipulation (FERM)**. FERM utilizes off-policy RL with
59 data augmentation along with unsupervised pre-
60 training to learn efficiently with a simple three-
61 staged procedure. First, a small number of (10)
62 demonstrations are collected and stored in a re-
63 play buffer. Second, the convolutional encoder
64 weights are initialized with unsupervised con-
65 trastive pre-training on the demonstration data.
66 Third, an off-policy RL algorithm is trained with
67 augmented images on both data collected online
68 during training and the initial demonstrations. Our
69 core contribution is the novel combination of contrastive pre-training, online data augmentations, and
70 utilizing a small number of demonstrations that together enable efficient real-robot learning from
71 pixels. In contrast, prior leading algorithms that utilize these components individually are unable to
72 learn efficiently on real robots.

73 We summarize our key contributions and benefits of the FERM algorithm: (1) *Data-efficiency*:
74 FERM enables learning optimal policies on 6 diverse manipulation tasks such as reaching, pushing,
75 moving, pulling a large object, flipping a switch, drawer opening **in 15-50 minutes of total training**
76 **time** for each task. (2) *Real-robot deployment*: FERM trains efficiently on real robotic hardware
77 while prior related approaches that were successful in simulation Laskin et al. (2020,?) fail to learn
78 robust real-robot policies. (3) *Simplicity*: FERM is a novel combination of existing ideas such as
79 contrastive pre-training, data augmentation, and demonstrations that results in a simple and easy to
80 reproduce algorithm. (4) *General & lightweight setup*: Our setup requires a robot, one GPU, two
81 cameras, a handful of demonstrations, and a sparse reward function. These requirements are quite
82 lightweight relative to setups that rely on Sim2Real, motion capture, multiple robots, or engineering
83 dense rewards. To the best of our knowledge, this work is the first to show how recent advances in
84 contrastive learning and data augmentation can enable efficient real-robot reinforcement learning
85 from pixels.

86 2 Background

87 **Soft Actor Critic**: The Soft Actor Critic (SAC) Haarnoja et al. (2018) is an off-policy RL algorithm
88 that jointly learns an action-conditioned state value function through Q learning and a stochastic
89 policy by maximizing expected returns. SAC is a state-of-the-art model-free RL algorithm for
90 continuous control from state Haarnoja et al. (2018) and, in the presence of data augmentations,
91 from pixels as well Laskin et al. (2020); Kostrikov et al. (2020). In simulated benchmarks, such as
92 DeepMind control Tassa et al. (2018), SAC is as data-efficient from pixels as it is from state. For this

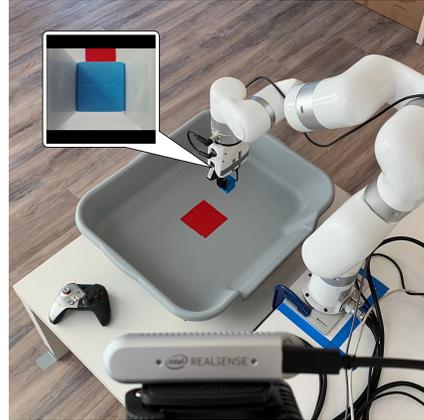


Figure 1: Framework for Efficient Robotic Manipulation (FERM) enables robotic agents to learn skills directly from pixels in less than one hour of training. Our setup requires a robotic arm, two cameras, and a joystick to provide 10 demonstrations.

93 reason, we utilize it as our base RL algorithm for sparse-reward manipulation in this work. As an
 94 actor-critic method, SAC learns an actor policy π_θ and an ensemble of critics Q_{ϕ_1} and Q_{ϕ_2} .

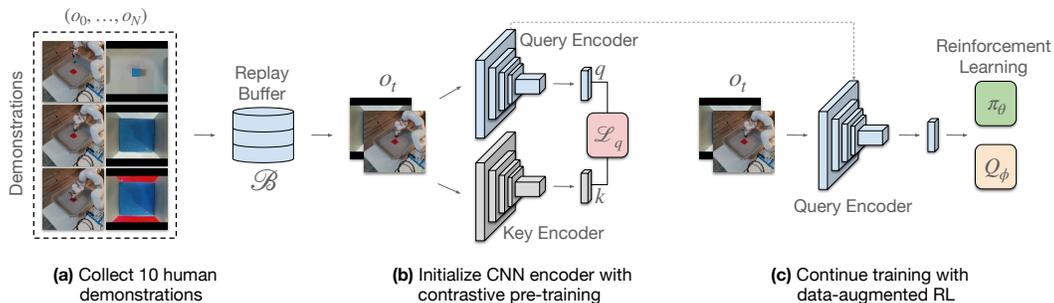


Figure 2: The FERM architecture. **(a)** Demonstrations are collected, and stored in a replay buffer. **(b)** The observations from the demonstrations are used to pre-train the encoder with a contrastive loss. **(c)** The encoder and replay buffer are then used to train an RL agent using an off-policy data-augmented RL algorithm.

95 To learn the actor policy, samples are collected stochastically from π_θ such that $a_\theta(o, \xi) \sim$
 96 $\tanh(\mu_\theta(o) + \sigma_\theta(o) \odot \xi)$, where $\xi \sim \mathcal{N}(0, I)$ is a sample from a normalized Gaussian noise vector,
 97 and then trained to maximize the expected return as shown in eq. 1.

$$\mathcal{L}(\theta) = \mathbb{E}_{a \sim \pi} [Q^\pi(o, a) - \alpha \log \pi_\theta(a|o)] \quad (1)$$

98 Simultaneously to learning the policy, SAC also trains the critics Q_{ϕ_1} and Q_{ϕ_2} to minimize the
 99 Bellman equation in Equation 2. Here, a transition $t = (o, a, o', r, d)$ is sampled from the replay
 100 buffer \mathcal{B} , where (o, o') are consecutive timestep observations, a is the action, r is the reward, and d is
 101 the terminal flag.

$$\mathcal{L}(\phi_i, \mathcal{B}) = \mathbb{E}_{t \sim \mathcal{B}} \left[(Q_{\phi_i}(o, a) - (r + \gamma(1-d)Q_{\text{targ}}))^2 \right] \quad (2)$$

102 The function Q_{targ} is the target value that the critics are trained to match, defined in Equation 3.
 103 The target is the entropy regularized exponential moving average (EMA) of the critic ensemble
 104 parameters, which we denote as \bar{Q}_ϕ .

$$Q_{\text{targ}} = \left(\min_{i=1,2} \bar{Q}_{\phi_i}(o', a') - \alpha \log \pi_\theta(a'|o') \right) \quad (3)$$

105 where (a', o') are the consecutive timestep action and observation, and α is a positive action-entropy
 106 coefficient. A non-zero action-entropy term improves exploration – the higher the value of α to more
 107 entropy maximization is prioritized over optimizing the value function.

108 **Unsupervised Contrastive Pretraining:** Contrastive learning Hadsell et al. (2006); LeCun et al.
 109 (2006); van den Oord et al. (2018); Wu et al. (2018); He et al. (2019); Chen et al. (2020); He
 110 et al. (2020); Hénaff et al. (2019) aims to maximize agreement between positive examples in
 111 data while minimizing agreement between negative examples. Contrastive methods require the
 112 specification of *query-key* pairs, also known as *anchors* and *positives*, which are similar data pairs
 113 whose agreement needs to be maximized. Given a query q and a key k , we seek to maximize the
 114 score $f_{\text{score}}(q, k)$ between them while minimizing them between the query q and negative examples in
 115 the dataset k_- . The score function is most often represented as an inner product, such as a dot product
 116 $f_{\text{score}}(q, k) = q^T k$ Wu et al. (2018); He et al. (2019) or a bilinear product $f_{\text{score}}(q, k) = q^T W k$
 117 van den Oord et al. (2018); Hénaff et al. (2019), while other Euclidean metrics are also available
 118 Schroff et al. (2015); Wang and Gupta (2015). Modern contrastive approaches Chen et al. (2020); He
 119 et al. (2020); Hénaff et al. (2019); Laskin et al. (2020) employ the InfoNCE loss van den Oord et al.

120 (2018), which is described in Equation 4 and can also be interpreted as a multi-class cross entropy
 121 classification loss with K classes.

$$\mathcal{L}_q = \log \frac{\exp(q^T W k)}{\exp\left(\sum_{i=0}^K \exp(q^T W k_i)\right)} \quad (4)$$

122 In the computer vision setting, a simple and natural choice of query-key specification is to define
 123 queries and keys as two data augmentations of the same image. This approach, called instance
 124 discrimination, is used in most of the state-of-the-art representation learning methods for static
 125 images Chen et al. (2020); He et al. (2020) as well as RL from pixels Laskin et al. (2020). In the
 126 minibatch setting, which we also employ in this work, the InfoNCE loss is computed by sampling
 127 $K = \{x_1, \dots, x_K\}$ images from the dataset, generating queries $Q = \{q_1, \dots, q_K\}$ and keys
 128 $K = \{k_1, \dots, k_K\}$ with stochastic data augmentations $q_i, k_i = \text{aug}(x_i)$, and using each augmented
 129 datapoint x_i as positives while the rest of the images are negatives.

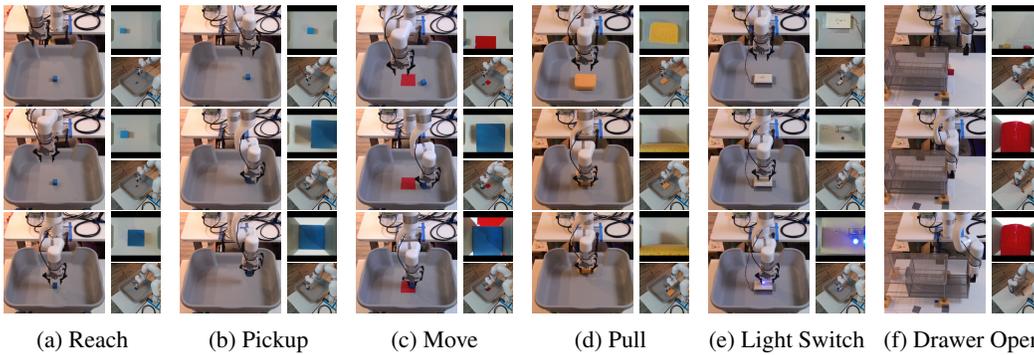


Figure 3: The set of real world tasks used in this work, along with their pixel observations. Each column shows initial, intermediate, and completion states of a rollout during evaluation of our optimal policy. The right two images comprise the processed camera image input, which are concatenated and used as the observational input for the RL agent. The sparse reward is only given when the robot completes the task. FERM is able to solve all 6 tasks within an hour, using only 10 demonstrations.

130 3 Method

131 Our proposed framework, shown
 132 in Figure 2, combines demonstra-
 133 tions, unsupervised pre-training,
 134 and off-policy model-free RL
 135 with data augmentation into one
 136 holistic Framework. FERM has
 137 three distinct steps – (i) minimal
 138 collection of demonstrations (ii)
 139 encoder initialization with unsu-
 140 pervised pre-training and (iii) on-
 141 line policy learning through RL
 142 with augmented data – which we
 143 describe in detail below.

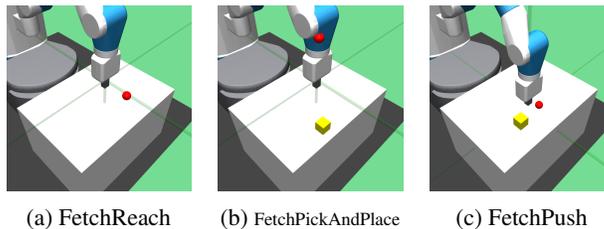


Figure 4: Simulated environments from OpenAI gym Brockman et al. (2016) used in addition to our real robot experiments. We use the Fetch Gym Suite to investigate the core components of FERM .

144 **Minimal Collection of Demonstrations:** We initialize the replay buffer with a small number of
 145 expert demonstrations (we found 10 to be sufficient) for each task. Demonstrations are collected with
 146 a joystick controller, shown in Figure 1. Our goal is to minimize the total time required to acquire a
 147 skill for an RL agent, including both policy training as well as time required to collect demonstrations.
 148 While collecting a larger number of demonstrations certainly improves training speed, we find 10
 149 demonstrations is already sufficient to learn skills quickly (see Fig. 7). For real world experiments,
 150 collecting 10 expert demonstrations can be done within 10 minutes which includes the time needed to
 151 reset the environment after every demonstration.

152 **Unsupervised Encoder Pre-training:** After initializing the replay buffer with 10 demonstrations,
 153 we pre-train the convolutional encoder with instance-based contrastive learning, using stochastic
 154 random crop Laskin et al. (2020) to generate query-key pairs. The key encoder is an exponentially
 155 moving average of the query encoder He et al. (2020), and the similarity measure between query-key
 156 pairs is the bi-linear inner product van den Oord et al. (2018) shown in Equation 4. Note that the
 157 bi-linear inner product is only used to pre-train the encoder. After pre-training, the weight matrix in
 158 the bi-linear measure is discarded.

159 **Reinforcement Learning with Augmented Data:** After pre-training the convolutional encoder on
 160 offline demonstration data, we train a SAC Haarnoja et al. (2018) agent with data augmentation
 161 Laskin et al. (2020) as the robot interacts with the environment. Since the replay buffer was initialized
 162 with demonstrations and SAC is an off-policy RL algorithm, during each minibatch update the
 163 agent receives a mix of demonstration observations and observations collected during training when
 164 performing gradient updates. The image augmentation used during training is random crop – the
 165 same augmentation used during contrastive pre-training.

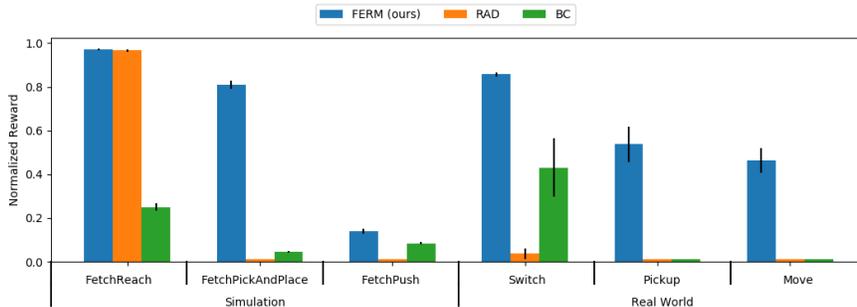


Figure 5: Baseline Comparisons. Shown are normalized rewards of the agent at the end of training for the simulated as well as the real robot results, as well as standard error. While FERM is able to learn all the tasks, the baseline RL agent (RAD) is unable to learn one sparse reward task without demonstrations. Conversely, with access to only 10 demonstrations, behavior cloning is unable to learn in the more difficult environments, and only succeeds on the simpler tasks (FetchReach, Light Switch). FERM and RAD are trained for 200k environment steps in simulated tasks, and until convergence for real world tasks (30 episodes for Switch and Pickup, 60 episodes for Move). BC is trained over the dataset for 200 epochs for both simulated and real world tasks. Simulated tasks are evaluated over 100 episodes, while real world tasks are evaluated over 30 episodes.

166 4 Experimental Evaluation

167 4.1 Experimental Setup

168 **Real robot:** We use the xArm xar robot for all real-world experiments. The end effector, a parallel
 169 two-jaw gripper, is position controlled with three degrees of freedom. At each step, the robot takes in
 170 an action containing the end effector and gripper aperture displacement.

171 **Operation space:** The range of motion of the gripper is confined to a 25 cm-high imaginary box
 172 above the manipulation surface. For majority of the tasks, objects are contained in a plastic tray
 173 approximately 40×34 cm in size measured at its bottom. Sponge padding was placed below the tray
 174 to absorb minor collisions between the gripper and the objects.

175 **Input:** We use two RGB cameras, one positioned over the shoulder for maximal view of the arm,
 176 and the other located within the gripper to provide a local object-level view. The over-the-shoulder
 177 camera is an Intel Realsense D415, with native resolution of 1280×720 . Specifically, we only utilize
 178 the RGB frames during both training as testing. Inside the gripper, we use an Arducam 8MP camera
 179 module configured to 640×480 in resolution. Image frames from both cameras are cropped and
 180 down-sampled to 100×100 pixels for use in our training algorithm.

181 **Demonstrations:** Using a Xbox controller xbo, we teleoperate the robot by supplying the end effector
 182 and gripper aperture displacement. Collecting demonstrations for each task requires less than **10**
 183 **minutes**, including resetting the environment.

184 **4.2 Environments and Baselines**

185 **Environments:** We evaluate FERM on six real-robotic manipulation tasks - reaching an object,
 186 picking up a block, moving a block to a target destination, pulling a large deformable object, flipping
 187 a switch, and opening a drawer. The block manipulation tasks (reach, pickup, move) are real-
 188 world adaptations of tasks from the OpenAI Gym Fetch suite Brockman et al. (2016). We utilize
 189 these three OpenAI gym environments for simulated environment experiments. Since our method
 190 uses demonstrations, we include pull, which has been used in prior work on imitation learning
 191 Rahmatizadeh et al. (2017); Florence et al. (2020). Flipping a switch is included as it demands
 192 precision, while drawer opening is a common task in existing simulated robotic benchmarks Yu et al.
 193 (2019). Details of task setup are provided in the supplementary material.

194 **Baselines:** We compare FERM to RAD Laskin et al. (2020), a leading supervised RL algorithm
 195 in simulated environments, and behavior cloning for our main results in Fig. 5. In the ablations
 196 section, we investigate each individual component of FERM . We investigate the contribution of
 197 each component of the FERM algorithm by removing one component - demonstrations, contrastive
 198 pre-training, or data augmentation - while keeping others fixed.

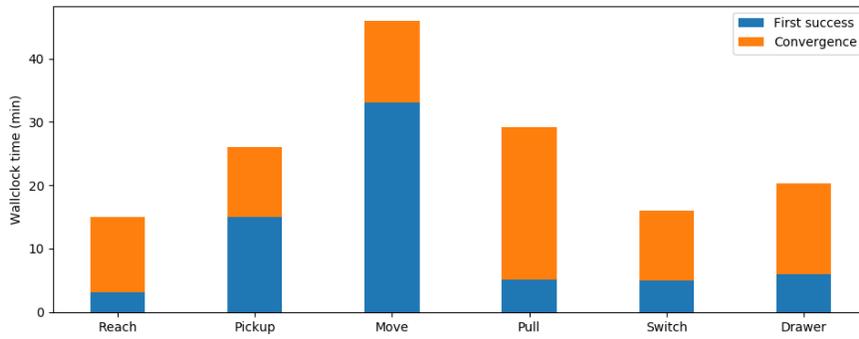


Figure 6: The speed at which our agents learn to complete the tasks. Plotted above are the times at which the policy first achieves a success, as well as when an optimal policy is learnt. Our method starts to complete the tasks in around 30 minutes of training, and as little as 3 minutes for simple tasks such as Reach.

Table 1: The success rates when evaluating the final policy learned by FERM over 30 episodes. Our method is able to achieve perfect success rate on the simpler tasks (Reach, Pickup, Light Switch, Drawer Open), and high success rates on the harder tasks (Move, Pull).

TASKS	REACH	PICKUP	MOVE	PULL	LIGHT SWITCH	DRAWER OPEN
# SUCCESSES (/30)	30	30	26	28	30	30
SUCCESS RATE (%)	100	100	86	93	100	100

199 **4.3 Results**

200 The main results of our investigation, including the time required to train an optimal policy as well the
 201 first successful task completion, are shown in Figure 5 and Table 1. We summarize the key findings
 202 below:

- 203 (i) On average, FERM enables a single robotic arm to learn optimal policies across all 6 tasks tested
 204 **within 30 minutes of training time** with a range of 15-50 minutes, which corresponds to to 20-80
 205 episodes of training. (see Fig. 6 and Table 1).
- 206 (ii) When evaluated on 3 simulated and 3 real-robot tasks, FERM substantially outperforms RAD and
 207 behavior cloning baselines (see Fig. 5).
- 208 (iii) The time to first successful task completion is **on average 11 minutes** with a range of 3-33
 209 minutes. The final policies achieve an **average success rate of 96%** with a range of 86-100% across
 210 the tasks tested, suggesting that they have converged to near-optimal solutions to the tasks.

211 (iv) Collecting demonstrations and contrastive pre-training does not introduce significant overhead in
 212 terms of time. Collecting 10 expert demonstrations with a joystick requires 10 minutes of human
 213 operation. Contrastive pre-training completes within one minute on a single NVIDIA 2080Ti GPU.
 214

215 (v) FERM solves all 6 tasks using the **same hyperparameters** and without altering the camera setup,
 216 which demonstrates the ease of use and generality of the framework.

217 Altogether, an RL agent trained with FERM is able to learn optimal policies for the 6 tasks ex-
 218 tremely efficiently. While prior work was able to solve dexterous manipulation tasks using RL with
 219 demonstrations in 2-3 hours of training Zhu et al. (2019), it also utilized dense rewards and more
 220 demonstrations. To the best of our knowledge, FERM is the first reinforcement learning method to
 221 solve a diverse set of sparse-reward robotic manipulation tasks directly from pixels in less than one
 222 hour.

223 4.4 Ablations

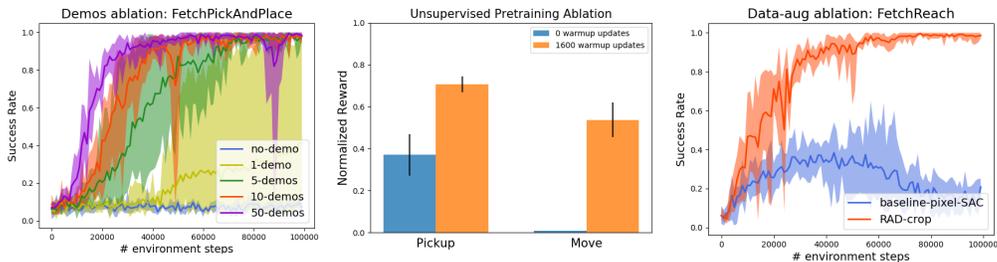


Figure 7: **Left:** We ablate the number of demonstrations required by FERM, and find that although the agent fails to learn with zero demonstrations, it can learn the PickAndPlace task efficiently using only 10 demonstrations. **Center:** We compare the performance of the move task with and without the use of pre-training on the real xArm robot. The plotted episode returns at convergence show that the contrastive pre-training substantially boosts performance. **Right:** Policy performance is measured by evaluation success rate. Using data augmentation, the agent achieves successful performance. Using non-augmented observations, the agent fails to learn the task.

224 In this section, we investigate how the three core components of FERM – demonstrations, contrastive
 225 pre-training, and data augmentation – contribute to the overall efficiency of the framework.

226 *How many demonstrations are needed?*

227 While sparse rewards are simpler to define, they pose an exploration challenge since the robot is
 228 unlikely to randomly stumble on a reward state. We address this issue by providing demonstrations to
 229 the RL agent. We ablate the number of demonstrations required to learn efficiently on the simulated
 230 pick and place task in Figure 7. We find that while the agent fails entirely with zero demonstrations,
 231 it is able to start learning the task with just one demonstration. While more demonstrations improve
 232 learning efficiency and reduce the variance of the policy, ten demonstrations suffice to learn quickly.
 233 We then evaluate the effectiveness of the 10 demonstrations by comparing our method to training
 234 behavior cloning. As shown in Figure 5, the 10 demonstrations are not enough to learn an effective
 235 policy. Refer to the supplementary material for further details.

236 *How important is unsupervised contrastive pre-training?* We next study the role of contrastive
 237 pre-training in FERM. We ablate our method with and without contrastive pre-training on the real
 238 world pickup and move task, shown in Figure 7, where we compare with (0), and without (1600)
 239 iterations of pre-training to initialize the encoder. With 1600 contrastive iterations, the agent is able
 240 to learn a successful policy while the other runs fail to learn. In the case of no pre-training at all, the
 241 agent is only able to succeed once during the entire hour of training.

242 *Is online data augmentation necessary?* To justify the use of data augmentation during online RL
 243 training, we compare the performance of SAC with and without data augmentation for a simple,
 244 dense reaching task. In the FetchReach environment, we use the dense reward $r = -d$ where
 245 d is the Euclidean distance between the gripper and the goal. As shown in Figure 7, without data

246 augmentation, the RL agent is unable to learn the simple task, and asymptotically collapses. This
247 motivates us to use data augmentation for our sparse reward tasks, which encounter even less learning
248 signal.

249 5 Related Work

250 **Imitation Learning:** Imitation learning is a framework for learning autonomous skills from demon-
251 strations. One of the simplest and perhaps most widely used forms of imitation learning is behavior
252 cloning (BC) where an agent learns a skill by regressing onto demonstration data. BC has been
253 successfully applied across diverse modalities including video games Ross et al. (2011), autonomous
254 navigation Pomerleau (1988); Bojarski et al. (2016), autonomous aviation Giusti et al. (2016), lo-
255 comotion Nakanishi et al. (2004); Kalakrishnan et al. (2009), and manipulation Duan et al. (2017);
256 Zhang et al. (2018); Young et al. (2020); Rahmatizadeh et al. (2017). Other imitation learning
257 approaches include Dataset Aggregation Ross et al. (2010), Inverse Reinforcement Learning Ng and
258 Russell (2000); Abbeel and Ng (2004), and Generative Adversarial Imitation Learning Ho and Ermon
259 (2016). A general limitation of imitation learning approaches is the requirement for a large number
260 of demonstrations for each task Sharma et al. (2018). Although recent advancements have shown that
261 imitation learning can learn with a much more modest amount of demonstrations Zhang et al. (2018);
262 Rahmatizadeh et al. (2017); Florence et al. (2020), FERM can learn in the same number of episodes,
263 of which the majority are spent with reinforcement learning.

264 **Reinforcement Learning:** Reinforcement Learning (RL) has been a promising approach for robotic
265 manipulation due to its ability to learn skills autonomously, but has not achieved widespread adoption
266 in real-world robotics. Recently, deep RL methods excelled at playing video games from pixels Mnih
267 et al. (2015); Berner et al. (2019) as well as learning robotic manipulation policies from visual input
268 Levine et al. (2015); Finn and Levine (2017); Haarnoja et al. (2018); Nair et al. (2018a). However,
269 widespread adoption of RL in real-world robotics has been bottle-necked due to the data-inefficiency
270 of the method, among other factors such as safety. Though there exist prior frameworks for efficient
271 position controlled robotic manipulation Zhu et al. (2019), they still require hours of training per task
272 and provide additional information such as a dense reward function. FERM is most closely related to
273 other methods that use RL with demonstrations. Prior methods Nair et al. (2018b); Rajeswaran et al.
274 (2017); Vecerík et al. (2017) solve robotic manipulation tasks from coordinate state input, rather than
275 image input, by initializing the replay buffer of an RL algorithm with demonstrations to overcome
276 the exploration problem in the sparse reward setting.

277 **Data Augmentation:** Image augmentation refers to stochastically altering images through transfor-
278 mations such as cropping, rotating, or color-jittering. It is widely used in computer vision architectures
279 including seminal works such as LeNet Lecun et al. (1998) and AlexNet Krizhevsky et al. (2017).
280 Data augmentation has played a crucial role in unsupervised representation learning in computer
281 vision Hénaff et al. (2019); He et al. (2020); Chen et al. (2020), while other works investigated
282 automatic generation of data augmentation strategies Cubuk et al. (2019). Data augmentation has
283 also been utilized in prior real robot RL methods Kalashnikov et al. (2018); however, the extent of its
284 significance for efficient training was not fully understood until recent works Laskin et al. (2020,?);
285 Kostrikov et al. (2020), which showed that carefully implemented data augmentation makes RL
286 policies from pixels as efficient as those from coordinate state. Finally, data augmentation has also
287 been shown to improve performance in imitation learning Young et al. (2020). In this work, data
288 augmentation comprises one of three components of a general framework for efficient learning.

289 **Unsupervised Representation Learning:** The goal of unsupervised representation learning is to
290 extract representations of high-dimensional unlabeled data that can then be used to learn downstream
291 tasks efficiently. Most relevant to our work is contrastive learning, which is a framework for learning
292 effective representations that satisfy similarity constraints between a pair of points in dataset. In
293 contrastive learning, latent embeddings are learned by minimizing the latent distance between similar
294 data points and maximizing them between dissimilar ones. Recently, a number of contrastive learning
295 methods Hénaff et al. (2019); He et al. (2019); Chen et al. (2020) have achieved state-of-the-art label-
296 efficient training in computer vision. A number of recent investigations in robotics have leveraged
297 contrastive losses to learn viewpoint invariant representations from videos Sermanet et al. (2018),
298 manipulate deformable objects Yan et al. (2020), and learn object representations Pirk et al. (2019).
299 In this work, we focus on instance-based contrastive learning Wu et al. (2018) similar to how it is

300 used in vision He et al. (2020); Chen et al. (2020) and RL on simulated benchmarks Laskin et al.
301 (2020); Stooke et al. (2020).

302 **6 Limitations**

303 Although FERM enables data-efficient deployment of RL onto real robots, the method also has a
304 number of limitations. First, like most RL algorithms, FERM may require assistance for resets, and
305 FERM policies can only solve the tasks that they were trained on and while they may display some
306 degree of generalization to small changes such as object shape or perturbations, we do not expect
307 FERM policies to generalize to qualitatively different tasks that were unseen during training. Second,
308 while the tasks considered in this paper are standard robotics evaluation tasks, they all have relatively
309 short horizons. Since FERM relies on a sparse reward signal to learn, we do not expect this framework
310 to succeed in long-horizon sparse reward tasks, where random interaction with the reward is unlikely.
311 Finally, we expect the performance of FERM to degrade if the visual conditions of the scene change
312 substantially, which is likely in non-lab settings with frequent background distractors and lighting
313 changes. Rather than addressing generalization to new tasks and visual settings or long-horizon
314 settings, this paper focuses on the data-efficiency problem of training RL policies on real robots.
315 We believe that data-efficient generalization and long-horizon problem solving are important open
316 problem in robot learning that we leave for future work.

317 **7 Conclusion and Future Work**

318 We present FERM , a framework that combines demonstrations, unsupervised learning, and RL, to
319 efficiently learn complex tasks in the real world. Using image input, our method is able to successfully
320 solve a diverse set of tasks, all using the same hyperparameters, and from sparse reward. Due to the
321 limited amount of supervision required, our work presents exciting avenues for applying RL to real
322 robots in a quick and efficient manner.

323 **References**

- 324 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Ried-
325 miller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement
326 learning. *Nature*, 518(7540):529, 2015.
- 327 C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi, Q. Fischer,
328 S. Hashme, C. Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint*
329 *arXiv:1912.06680*, 2019.
- 330 O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell,
331 T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai,
332 J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden,
333 Y. Sulsky, J. Molloy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama,
334 D. Wünsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu, D. Hassabis,
335 C. Apps, and D. Silver. Grandmaster level in StarCraft II using multi-agent reinforcement
336 learning. *Nature*, 575(7782):350–354, 2019. doi: 10.1038/s41586-019-1724-z. URL <https://doi.org/10.1038/s41586-019-1724-z>.
- 338 A. P. Badia, B. Piot, S. Kapturowski, P. Sprechmann, A. Vitvitskyi, D. Guo, and C. Blundell. Agent57:
339 Outperforming the atari human benchmark. In *International Conference on Machine Learning*,
340 2020.
- 341 D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser,
342 I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner,
343 I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the
344 game of Go with deep neural networks and tree search. *Nature*, 529:484–503, 2016. URL <http://www.nature.com/nature/journal/v529/n7587/full/nature16961.html>.
- 346 D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker,
347 M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. Driessche, T. Graepel, and D. Hassabis.
348 Mastering the game of go without human knowledge. *Nature*, 550:354–359, 10 2017. doi:
349 10.1038/nature24270.
- 350 J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart,
351 D. Hassabis, T. Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned
352 model. *arXiv preprint arXiv:1911.08265*, 2019.
- 353 J. Schulman, S. Levine, P. Abbeel, M. I. Jordan, and P. Moritz. Trust region policy optimization. In
354 *ICML*, pages 1889–1897, 2015.
- 355 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
356 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 357 M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas. Reinforcement learning with
358 augmented data. *arXiv preprint arXiv:2004.14990*, 2020.
- 359 D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
360 imagination. In *International Conference on Learning Representations*, 2020.
- 361 S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-End Training of Deep Visuomotor Policies.
362 *CoRR*, abs/1504.00702, 2015. URL <http://arxiv.org/abs/1504.00702>.
- 363 D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan,
364 V. Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic
365 manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- 366 L. Pinto and A. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot
367 hours. *ICRA*, 2016.
- 368 J. Mahler, F. T. Pokorny, B. Hou, M. Roderick, M. Laskey, M. Aubry, K. Kohlhoff, T. Kröger,
369 J. Kuffner, and K. Goldberg. Dex-net 1.0: A cloud-based network of 3d objects for robust grasp
370 planning using a multi-armed bandit model with correlated rewards. In *ICRA*, 2016.

- 371 S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic
372 grasping with deep learning and large-scale data collection. *ISER*, 2016.
- 373 A. Gupta, A. Murali, D. P. Gandhi, and L. Pinto. Robot learning in homes: Improving generalization
374 and reducing dataset bias. In *Advances in Neural Information Processing Systems*, pages 9094–
375 9104, 2018.
- 376 M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron,
377 M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint*
378 *arXiv:1808.00177*, 2018.
- 379 M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot,
380 M. Azar, and D. Silver. Rainbow: Combining improvements in deep reinforcement learning, 2017.
- 381 T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous
382 control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- 383 Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel,
384 A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- 385 M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,
386 O. Pieter Abbeel, and W. Zaremba. Hindsight experience replay. In *NeurIPS*, 2017.
- 387 T. Zhang, Z. McCarthy, O. Jow, D. Lee, K. Goldberg, and P. Abbeel. Deep imitation learning for
388 complex manipulation tasks from virtual reality teleoperation. *2018 IEEE International Conference*
389 *on Robotics and Automation (ICRA)*, pages 1–8, 2018.
- 390 J. Ho and S. Ermon. Generative Adversarial Imitation Learning. In D. D. Lee, M. Sugiyama, U. V.
391 Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*,
392 pages 4565–4573. Curran Associates, Inc., 2016. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/6391-generative-adversarial-imitation-learning.pdf)
393 [6391-generative-adversarial-imitation-learning.pdf](http://papers.nips.cc/paper/6391-generative-adversarial-imitation-learning.pdf).
- 394 Y. Duan, M. Andrychowicz, B. C. Stadie, J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba.
395 One-shot imitation learning, 2017.
- 396 C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine. One-Shot Visual Imitation Learning via Meta-
397 Learning. volume 78 of *Proceedings of Machine Learning Research*, pages 357–368. PMLR,
398 13–15 Nov 2017. URL <http://proceedings.mlr.press/v78/finn17a.html>.
- 399 S. Young, D. Gandhi, S. Tulsiani, A. Gupta, P. Abbeel, and L. Pinto. Visual Imitation Made Easy.
400 *CoRR*, abs/2008.04899, 2020. URL <https://arxiv.org/abs/2008.04899>.
- 401 M. Laskin, A. Srinivas, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforce-
402 ment learning. In *International Conference on Machine Learning*, 2020.
- 403 A. Stooke, K. Lee, P. Abbeel, and M. Laskin. Decoupling representation learning from reinforcement
404 learning, 2020.
- 405 I. Kostrikov, D. Yarats, and R. Fergus. Image augmentation is all you need: Regularizing deep
406 reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.
- 407 M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An
408 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,
409 2013.
- 410 T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta,
411 P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*,
412 2018.
- 413 R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In
414 *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*,
415 volume 2, pages 1735–1742. IEEE, 2006.
- 416 Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning.
417 2006.

- 418 A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding.
419 *arXiv preprint arXiv:1807.03748*, 2018.
- 420 Z. Wu, Y. Xiong, S. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance-level
421 discrimination. *arXiv preprint arXiv:1805.01978*, 2018.
- 422 K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual
423 representation learning. *arXiv preprint arXiv:1911.05722*, 2019.
- 424 T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of
425 visual representations. In *International Conference on Machine Learning*, 2020.
- 426 K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual
427 representation learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
428 2020.
- 429 O. J. Hénaff, A. Srinivas, J. De Fauw, A. Razavi, C. Doersch, S. Eslami, and A. v. d. Oord. Data-
430 efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*,
431 2019.
- 432 F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and
433 clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
434 pages 815–823, 2015.
- 435 X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015.
- 436 G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai
437 gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 438 xArm 7. URL <https://store.ufactory.cc/products/xarm-7-2020>.
- 439 Xbox wireless controller. URL [https://www.xbox.com/en-US/accessories/
440 controllers/xbox-wireless-controller](https://www.xbox.com/en-US/accessories/controllers/xbox-wireless-controller).
- 441 R. Rahmatizadeh, P. Abolghasemi, L. Bölöni, and S. Levine. Vision-Based Multi-Task Manipulation
442 for Inexpensive Robots Using End-To-End Learning from Demonstration. *CoRR*, abs/1707.02920,
443 2017. URL <http://arxiv.org/abs/1707.02920>.
- 444 P. R. Florence, L. Manuelli, and R. Tedrake. Self-Supervised Correspondence in Visuomotor Policy
445 Learning. *IEEE Robotics Autom. Lett.*, 5(2):492–499, 2020. doi: 10.1109/LRA.2019.2956365.
446 URL <https://doi.org/10.1109/LRA.2019.2956365>.
- 447 T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-World: A Benchmark
448 and Evaluation for Multi-Task and Meta Reinforcement Learning. In *Conference on Robot
449 Learning (CoRL)*, 2019. URL <https://arxiv.org/abs/1910.10897>.
- 450 H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar. Dexterous manipulation with deep
451 reinforcement learning: Efficient, general, and low-cost. In *2019 International Conference on
452 Robotics and Automation (ICRA)*, pages 3651–3657, 2019. doi: 10.1109/ICRA.2019.8794102.
- 453 S. Ross, G. J. Gordon, and D. Bagnell. A Reduction of Imitation Learning and Structured Prediction
454 to No-Regret Online Learning. In G. J. Gordon, D. B. Dunson, and M. Dudík, editors, *Proceedings
455 of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011,
456 Fort Lauderdale, USA, April 11-13, 2011*, volume 15 of *JMLR Proceedings*, pages 627–635.
457 [JMLR.org](http://proceedings.mlr.press/v15/ross11a/ross11a.pdf), 2011. URL [http://proceedings.mlr.press/v15/ross11a/
458 pdf](http://proceedings.mlr.press/v15/ross11a/ross11a.pdf).
- 459 D. Pomerleau. ALVINN: an autonomous land vehicle in a neural network. In D. S. Touretzky, editor,
460 *Advances in Neural Information Processing Systems 1, [NIPS Conference, Denver, Colorado, USA,
461 1988]*, pages 305–313. Morgan Kaufmann, 1988. URL [http://papers.nips.cc/paper/
462 95-alvinn-an-autonomous-land-vehicle-in-a-neural-network](http://papers.nips.cc/paper/95-alvinn-an-autonomous-land-vehicle-in-a-neural-network).
- 463 M. Bojarski, D. D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort,
464 U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba. End to End Learning for Self-Driving Cars.
465 *CoRR*, abs/1604.07316, 2016. URL <http://arxiv.org/abs/1604.07316>.

- 466 A. Giusti, J. Guzzi, D. C. Ciresan, F. He, J. P. Rodriguez, F. Fontana, M. Faessler, C. Forster,
467 J. Schmidhuber, G. D. Caro, D. Scaramuzza, and L. M. Gambardella. A Machine Learning
468 Approach to Visual Perception of Forest Trails for Mobile Robots. *IEEE Robotics Autom. Lett.*, 1
469 (2):661–667, 2016. doi: 10.1109/LRA.2015.2509024. URL [https://doi.org/10.1109/
470 LRA.2015.2509024](https://doi.org/10.1109/LRA.2015.2509024).
- 471 J. Nakanishi, J. Morimoto, G. Endo, G. Cheng, S. Schaal, and M. Kawato. Learning from demon-
472 stration and adaptation of biped locomotion. *Robotics Auton. Syst.*, 47(2-3):79–91, 2004. doi: 10.
473 1016/j.robot.2004.03.003. URL <https://doi.org/10.1016/j.robot.2004.03.003>.
- 474 M. Kalakrishnan, J. Buchli, P. Pastor, and S. Schaal. Learning locomotion over rough terrain using
475 terrain templates. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems,
476 October 11-15, 2009, St. Louis, MO, USA*, pages 167–172. IEEE, 2009. doi: 10.1109/IROS.2009.
477 5354701. URL <https://doi.org/10.1109/IROS.2009.5354701>.
- 478 S. Ross, G. Gordan, and A. Bagnell. A reduction of imitation learning and structured prediction to
479 no-regret online learning. In *arXiv preprint arXiv:1011.0686*, 2010.
- 480 A. Y. Ng and S. J. Russell. Algorithms for inverse reinforcement learning. In P. Langley, editor,
481 *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000),
482 Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pages 663–670. Morgan Kaufmann,
483 2000.
- 484 P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In C. E. Brodley,
485 editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004),
486 Banff, Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding
487 Series*. ACM, 2004. doi: 10.1145/1015330.1015430. URL [https://doi.org/10.1145/
488 1015330.1015430](https://doi.org/10.1145/1015330.1015430).
- 489 P. Sharma, L. Mohan, L. Pinto, and A. Gupta. Multiple interactions made easy (mime): Large scale
490 demonstrations data for imitation. *arXiv preprint arXiv:1810.07121*, 2018.
- 491 C. Finn and S. Levine. Deep visual foresight for planning robot motion. In *2017 IEEE International
492 Conference on Robotics and Automation (ICRA)*, pages 2786–2793. IEEE, 2017.
- 493 A. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine. Visual Reinforcement Learning with
494 Imagined Goals. In S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi,
495 and R. Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Confer-
496 ence on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018,
497 Montréal, Canada*, pages 9209–9220, 2018a. URL [http://papers.nips.cc/paper/
498 8132-visual-reinforcement-learning-with-imagined-goals](http://papers.nips.cc/paper/8132-visual-reinforcement-learning-with-imagined-goals).
- 499 A. Nair, B. McGrew, M. Andrychowicz, W. Zaremba, and P. Abbeel. Overcoming Exploration
500 in Reinforcement Learning with Demonstrations. In *2018 IEEE International Conference on
501 Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, pages 6292–6299.
502 IEEE, 2018b. doi: 10.1109/ICRA.2018.8463162. URL [https://doi.org/10.1109/ICRA.
503 2018.8463162](https://doi.org/10.1109/ICRA.2018.8463162).
- 504 A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning
505 complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv
506 preprint arXiv:1709.10087*, 2017.
- 507 M. Vecerík, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and
508 M. A. Riedmiller. Leveraging Demonstrations for Deep Reinforcement Learning on Robotics
509 Problems with Sparse Rewards. *CoRR*, abs/1707.08817, 2017. URL [http://arxiv.org/
510 abs/1707.08817](http://arxiv.org/abs/1707.08817).
- 511 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
512 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
- 513 A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional
514 neural networks. *Commun. ACM*, 60(6):84–90, 2017. doi: 10.1145/3065386. URL [http://
515 //doi.acm.org/10.1145/3065386](http://doi.acm.org/10.1145/3065386).

- 516 E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. AutoAugment: Learning Augmentation
517 Strategies From Data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*
518 *2019, Long Beach, CA, USA, June 16-20, 2019*, pages 113–123. Computer Vision Foundation /
519 IEEE, 2019. doi: 10.1109/CVPR.2019.00020. URL [http://openaccess.thecvf.com/
520 content_CVPR_2019/html/Cubuk_AutoAugment_Learning_Augmentation_
521 Strategies_From_Data_CVPR_2019_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Cubuk_AutoAugment_Learning_Augmentation_Strategies_From_Data_CVPR_2019_paper.html).
- 522 P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine. Time-Contrastive
523 Networks: Self-Supervised Learning from Video. In *2018 IEEE International Conference on*
524 *Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, pages 1134–1141.
525 IEEE, 2018. doi: 10.1109/ICRA.2018.8462891. URL [https://doi.org/10.1109/ICRA.
526 2018.8462891](https://doi.org/10.1109/ICRA.2018.8462891).
- 527 W. Yan, A. Vangipuram, P. Abbeel, and L. Pinto. Learning Predictive Representations for Deformable
528 Objects Using Contrastive Estimation. *CoRR*, abs/2003.05436, 2020. URL [https://arxiv.
529 org/abs/2003.05436](https://arxiv.org/abs/2003.05436).
- 530 S. Pirk, M. Khansari, Y. Bai, C. Lynch, and P. Sermanet. Online Object Representations with
531 Contrastive Learning. *CoRR*, abs/1906.04312, 2019. URL [http://arxiv.org/abs/1906.
532 04312](http://arxiv.org/abs/1906.04312).

533 Checklist

- 534 1. For all authors...
- 535 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
536 contributions and scope? [Yes]
- 537 (b) Did you describe the limitations of your work? [Yes]
- 538 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 539 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
540 them? [Yes]
- 541 2. If you are including theoretical results...
- 542 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 543 (b) Did you include complete proofs of all theoretical results? [N/A]
- 544 3. If you ran experiments...
- 545 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
546 mental results (either in the supplemental material or as a URL)? [Yes]
- 547 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
548 were chosen)? [Yes]
- 549 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
550 ments multiple times)? [Yes]
- 551 (d) Did you include the total amount of compute and the type of resources used (e.g., type
552 of GPUs, internal cluster, or cloud provider)? [Yes]
- 553 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 554 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 555 (b) Did you mention the license of the assets? [Yes]
- 556 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 557 (d) Did you discuss whether and how consent was obtained from people whose data you’re
558 using/curating? [N/A]
- 559 (e) Did you discuss whether the data you are using/curating contains personally identifiable
560 information or offensive content? [N/A]
- 561 5. If you used crowdsourcing or conducted research with human subjects...
- 562 (a) Did you include the full text of instructions given to participants and screenshots, if
563 applicable? [N/A]

564
565
566
567

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]