

Learning Robot Manipulation from Cross-Morphology Demonstration

Anonymous Author(s)

Affiliation

Address

email

Abstract: Some Learning from Demonstrations (LfD) methods handle small mismatches in the action spaces of the teacher and student. Here we address the case where the teacher’s morphology is substantially different from that of the student. Our framework, Morphological Adaptation in Imitation Learning (**MAIL**), bridges this gap allowing us to train an agent from demonstrations by other agents with significantly different morphologies. **MAIL** learns from suboptimal demonstrations, so long as they provide *some* guidance towards a desired solution. We demonstrate **MAIL** on manipulation tasks with rigid and deformable objects including 3D cloth manipulation interacting with rigid obstacles. We train a visual control policy for a robot with one end-effector using demonstrations from a simulated agent with two end-effectors. **MAIL** shows up to 24% improvement in a normalized performance metric over LfD and non-LfD baselines. It is deployed to a real Franka Panda robot, handles multiple variations in properties for objects (size, rotation, translation), and cloth-specific properties (color, thickness, size, material). We show generalizability to morphology adaptation from n -to- m end-effectors, in a rearrangement task executed in simulation and the real world. An overview is on this [website](#).

Keywords: Imitation from Observation, Learning from Demonstration

1 Introduction

Learning from Demonstration (LfD) [1, 2] is a set of supervised learning methods where a teacher (often, but not always, a human) demonstrates a task, and a student (usually a robot) uses this information to learn to perform the same task. Some LfD methods cope with small morphological mismatches between the teacher and student [3, 4] (*e.g.*, five-fingered hand to two-fingered gripper). However, they typically fail for a large mismatch (*e.g.*, bimanual human demonstration to a robot arm with one gripper). The key difference is that to reproduce the transition from a demonstration state to the next, no single student action suffices - a sequence of actions may be needed.

Supervised methods are appealing where demonstration-free methods [5] do not converge or underperform [6] and purely analytical approaches are computationally infeasible [7, 8]. In such settings, human demonstrations of complex tasks are often readily available *e.g.*, it is straightforward for a human to show a robot how to fold a cloth. An LfD-based imitation learning approach is appealing in such settings *provided* we allow the human demonstrator to use their body in the way they find most convenient (*e.g.*, using two hands to hang a cloth on a clothesline to dry). This requirement induces a potentially large morphology mismatch - we want to learn and execute complex tasks with deformable objects on a single manipulator robot using natural human demonstrations.

We propose a framework, Morphological Adaptation in Imitation Learning (**MAIL**), to bridge this mismatch. **MAIL** enables policy learning for a robot with m end-effectors from teachers with n end-effectors. It does not require demonstrator actions, only the states of the objects in the environment

making it potentially useful for a variety of end-effectors (pickers, suction gripper, two-fingered grippers, or even hands). It uses trajectory optimization to convert state-based demonstrations into (suboptimal) trajectories in the student’s morphology. The optimization uses a learned (forward) dynamics model to trade accuracy for speed, especially useful for tasks with high-dimensional state and observation spaces. The trajectories are then used by an LfD method, which is adapted to work with suboptimal demonstrations and improve upon them by interacting with the environment.

Though the original demonstrations contain states, we generalize the solution to work with image observations in the final policy. We showcase our method on challenging cloth manipulation tasks (Sec. 4.1) for a robot with one end-effector, using image observations, shown in Fig. 1. This setting is challenging for multiple reasons. First, cloth manipulation is easy for bimanual human demonstrators but challenging for a one-handed agent (even humans find cloth manipulation non-trivial with one hand). Second, deformable objects exist in a continuous state space; image observations in this setting are also high-dimensional. Third, the cloth being manipulated makes a large number of contacts (hundreds) that are made/broken per time step. These can significantly slow down simulation, and consequently learning and optimization. We make the following contributions:

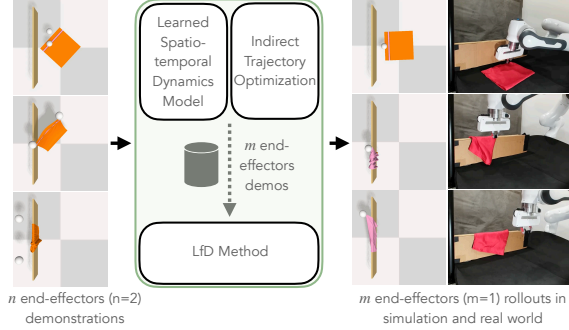


Figure 1: **MAIL** generalizes LfD to large morphological mismatches between teacher and student in difficult manipulation tasks. We show an example task: hang a cloth to dry on a plank (DRY CLOTH). The demonstrations are bimanual, yet the robot learns to execute the task with a single arm and gripper. The learned policy transfers to the real world and is robust to object variations.

1. We propose a novel framework, **MAIL**, that bridges the large morphological mismatch in LfD. **MAIL** trains a robot with m end-effectors to learn manipulation from demonstrations with a different (n) number of end-effectors (n -to- m end-effector transfer).
2. We demonstrate **MAIL** on challenging cloth manipulation tasks on a robot with one end-effector. Our tasks have a high-dimensional (> 15000) state space, with several 100 contacts being made/broken per step, and are non-trivial to solve with one end-effector. Our learned agent outperforms baselines by up to 24% on a normalized performance metric and transfers zero-shot to a real robot. We introduce a new variant of 3D cloth manipulation with obstacles - DRY CLOTH.
3. We illustrate how **MAIL** can handle general instances of n -to- m end-effector transfer, with a simple rearrangement task with three rigid bodies, in simulation and the real world. This task illustrates a 3-to-2, 3-to-1, and 2-to-1 end-effector transfer.

2 Related Work

Imitation Learning and Reinforcement Learning with Demonstrations: Imitation learning methods [9, 10, 11, 12, 13] and methods that combine reinforcement learning and demonstrations [14, 15, 1, 2] have shown excellent results in learning a mapping between observations and actions from demonstrations. However, their objective function requires access to the demonstrator’s ground truth actions for optimization. This is infeasible for cross-morphology transfer due to action space mismatch. To work around this, prior works have proposed systems for teachers to provide demonstrations in the students’ morphology [16] which limits the ability of teachers to efficiently provide data. Similar to imitation learning, offline RL [17, 18, 19] learns from demonstrations stored in a dataset without online environment interactions. While offline RL can work with large datasets of diverse rollouts to produce generalizable policies [20, 21], it requires the availability of rollouts that have the same action space as the learning agent. **MAIL** learns across morphologies and is not affected by this limitation.

Imitation from Observation: Imitation from observation (IfO) methods [3, 9, 22, 23, 24, 25] learn from the states of the demonstration; they do not use state-action pairs. In [26], an approach is proposed to learn repetitive actions using Dynamic Movement Primitives [27] and Bayesian optimization to maximize the similarity between human demonstrations and robot actions. Many IfO methods [3, 23, 24, 28] assume that the student can take a single action to transition from the demonstration’s current state to the next state. Some methods [3, 23] use this to train an inverse dynamics model to infer actions. Others extract keypoints from the observations and compute actions by subtracting consecutive keypoint vectors. However, when the student has a different space than the teacher, it may require more than one action for the student to reach consecutive demonstration states. For example, in an object rearrangement task, a two-picker teacher agent can move two objects with one pick-place action. But a one-picker student will need two or more actions to achieve the same result. Zero-shot visual imitation [9] assumes that the statistics of visual observations and agents observations will be similar. However, when solving a task with different numbers of arms, some intermediate states will not be seen in teacher demonstrations. State-of-the-art learning from observation methods [25, 29] have made significant advancements in exploiting information between states. However, their tasks have much longer horizons, hence more states and learning signals than ours. Whether these methods work well on short-horizon, difficult manipulation tasks is uncertain. To address this and provide a meaningful comparison, we conducted experiments to compare **MAIL** with these methods (Sec. 4).

Trajectory Optimization: Trajectory optimization algorithms [30, 8, 31] optimize a trajectory by minimizing a cost function, subject to a set of constraints. It has been used for manipulation of rigid and deformable objects [7], even through contact [32] using complementarity constraints [33]. Indirect trajectory optimization only optimizes the actions of a trajectory and uses a simulator for the dynamics instead of adding dynamics constraints at every step.

Learned Dynamics: Learning dynamics models is useful when there is no simulator, or if the simulator is too slow or too inaccurate. Learned models have been used with MPC to speed up prediction times [34, 35, 36]. A common use case is model-based RL [37], where learning the dynamics is part of the algorithm and has been shown to learn dynamics from states and pixels [38] and applied to real-world tasks [39].

3 Formulation and Approach

3.1 Preliminaries

We formulate the problem as a POMDP with state $\mathbf{s} \in \mathcal{S}$, action $\mathbf{a} \in \mathcal{A}$, observation $\mathbf{o} \in \mathcal{O}$, transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$, horizon H , discount factor γ and reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The discounted return at time t is $R_t = \sum_{i=t}^H \gamma^i r(\mathbf{s}_i, \mathbf{a}_i)$ and $\mathbf{s}_i \sim \mathcal{T}(\mathbf{s}_{i-1}, \mathbf{a}_{i-1})$. A task is instantiated with a variant sampled from the task distribution, $\mathbf{v} \sim \mathcal{V}$. The initial environment state depends on the task variant, $\mathbf{s}_0(\mathbf{v}), \mathbf{v} \sim \mathcal{V}$. We train a policy π_θ to maximize expected reward $J(\pi_\theta)$ of an episode over task variants \mathbf{v} , $J(\pi_\theta) = \mathbb{E}_{\mathbf{v} \sim \mathcal{V}}[R_0]$, subject to initial state $\mathbf{s}_0(\mathbf{v})$ and the dynamics from \mathcal{T} . For a method overview see Fig. 2.

For an agent with morphology M , we differentiate between datasets available as demonstrations (\mathcal{D}_{Demo}^M) and those that are optimized (\mathcal{D}_{Optim}^M). For our cloth environments, our teacher morphology is two-pickers ($M = 2p$) and student morphology is one-picker ($M = 1p$). We assume the demonstrations are from teachers with a morphology that can be different from the student (and from each other). We refer to these as *teacher* demonstrations, $\mathcal{D}_{Teacher}$, to emphasize that they do not necessarily come from an expert or an oracle. Further, these can be suboptimal. The demonstrations are state trajectories $\tau_T = (\mathbf{s}_0, \dots, \mathbf{s}_{H-1})$. The teacher dataset is made up of K_T such trajectories, $\mathcal{D}_{Teacher} = \{\tau_{T,i}\} \forall i = 1, \dots, K_T$, using a few task variations from the task distribution $\mathbf{v}_d \sim \mathcal{V}$.

We now discuss the components of **MAIL**. The user provides teacher demonstrations $\mathcal{D}_{Teacher}$. First, we create a dataset of random actions, \mathcal{D}_{Random} , and use it to train a dynamics model, \mathcal{T}_ψ . The learned dynamics are not task-specific and depend on the objects in the environment. \mathcal{T}_ψ reduces

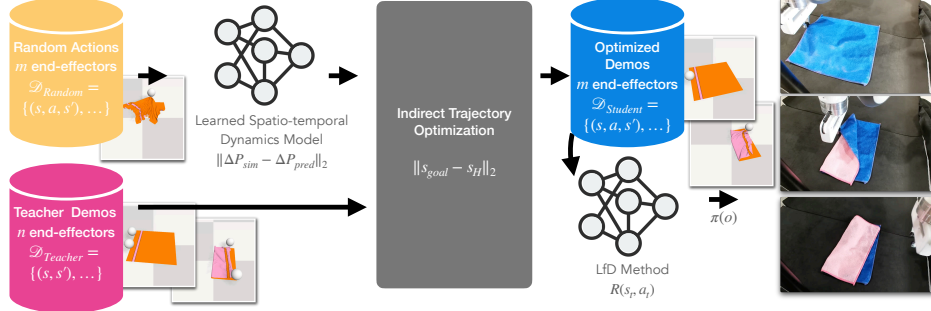


Figure 2: An example cloth folding task with demonstrations from a teacher with $n = 2$ end-effectors, deployed on a Franka Panda with $m = 1$ end-effector (parallel-jaw gripper). We train a network to predict the forward dynamics of the object being manipulated in simulation, using a random action dataset \mathcal{D}_{Random} . Given teacher demonstrations we use indirect trajectory optimization to find student actions that solve the task. Finally, we pass the optimized dataset $\mathcal{D}_{Student}$ to a downstream LfD method that combines imitation and RL to get a final policy π that generalizes to task variations and extends task learning to image space, enabling real-world deployment.

computational cost when dealing with contact-rich simulations like cloth manipulation (Sec. 4.1). Next, we convert each teacher demonstration to a trajectory suitable for the student’s morphology. For our tasks, we find gradient-free indirect trajectory optimization [31] performs the best (Appendix Sec. A.2.1). We used \mathcal{T}_ψ for this optimization as it provides the appropriate speed-accuracy trade-off. The optimization objective is to match with object states in the demonstration (we cannot match demonstration actions across morphologies). We combine these optimized trajectories to create a dataset $\mathcal{D}_{Student}$ for the student. Finally, we pass $\mathcal{D}_{Student}$ to a downstream LfD method to learn a policy π that generalizes from the task variations in $\mathcal{D}_{Teacher}$ to the task distribution \mathcal{V} . It also extends π to use image observations and deploys zero-shot on a real robot (rollouts in Fig. 5).

3.2 Learned Spatio-temporal Dynamics Model

MAIL uses trajectory optimization to convert demonstrations into (suboptimal) trajectories in the student’s morphology. This can be prohibitively slow for large state spaces and complex tasks such as cloth manipulation. Robotic simulators have come a long way in advancing fidelity and speed, but simulating complex deformable objects and contact-rich manipulation still requires significant computation making optimization intractable for challenging simulations. We use the NVIDIA FLeX simulator that is based on extended position-based dynamics [40]. We learn a CNN-LSTM based spatio-temporal forward dynamics model with parameters ψ , \mathcal{T}_ψ , to approximate cloth dynamics, \mathcal{T} . This offers a speed-accuracy trade-off with a tractable computation time in environments with large state spaces and complex dynamics. The states of objects are represented as N particle positions: $\mathbf{s} = P = \{p_i\}_{i=1\dots N}$. Each particle state consists of its x, y, and z coordinates. For each task, we generate a corpus of random pick-and-place actions and store them in the dataset $\mathcal{D}_{Random} = \{d_i\}$, where $i = 1, \dots, K_R$ and $d_i = (P_i, a_i, P'_i)$. For each datum i , we feed P_i to the CNN network to extract features of particle connectivity. These features are concatenated with a_i and input to the LSTM model to extract features based on the previous particle positions. A fully connected layer followed by layer normalization and \tanh activation is used to learn the non-linear combinations of features. The outputs are the predicted particle displacements. The objective function is the distance between predicted and ground-truth particle displacements, $\|\Delta P_{sim} - \Delta P_{pred}\|_2$. Here $\Delta P_{sim} = \{\Delta p_i\}_{i=1,\dots,N}$ is obtained from the simulator and $\Delta p_i = p_{i+1} - p_i$ for every particle i .

Due to its simplicity, the CNN-LSTM dynamics model provides fast inference, compared to a simulator which may have to perform many collision checks at any time step. This speedup is crucial when optimizing over a large state space, as long as the errors in particle positions are tolerable. In our experiments, we were able to get 162 fps with \mathcal{T}_ψ , compared to 3.4 fps with the FLeX simulator (50x speed up) (Fig. 8). However, this stage is optional if the environment is low-dimensional, or if the simulation speed-up from inference is not significant. Simulation accuracy is important

when training a final policy, to provide *accurate* end-effector locations for execution on a real robot. Hence, the learned dynamics model is not used for training in the downstream LfD method.

3.3 Indirect Trajectory Optimization

We use indirect trajectory optimization [31] to find the open-loop action trajectory to match the teacher state trajectory, τ_T . This optimizes for the student’s actions while propagating the state with a simulator. We use the learned dynamics \mathcal{T}_ψ to give us fast, approximate optimized trajectories. This is in contrast to direct trajectory optimization (or collocation) that optimizes both states and actions at every time step. Direct trajectory optimization requires dynamics constraints to ensure consistency among states being optimized, which can be challenging for discontinuous dynamics. We use the Cross-Entropy Method (CEM) for optimization, and compare this against other methods, such as SAC (Appendix A.2.1). The optimization objective is to match the object’s goal state s_{goal} in the demonstration with the same task variant v_d . Formally, the problem is defined as:

$$\min_{\mathbf{a}_t} \|\mathbf{s}_{goal} - \mathbf{s}_H\|_2 \text{ subject to } \mathbf{s}_0 = \mathbf{s}_0(v_d) \text{ and } \mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t) \forall t = 0, \dots, H-1 \quad (1)$$

where \mathbf{s}_H is the predicted final state. Note that if τ_T has a longer time horizon, it would help to match intermediate states and use multiple-shooting methods. After optimizing the action trajectories for each demonstration $\tau_{T,i} \in \mathcal{D}_{Teacher}$, we use them with the simulator to obtain the optimized trajectories in the student’s morphology. These are combined to create the student dataset, $\mathcal{D}_{Student} = \{\tau_1, \tau_2, \tau_3, \dots\}$, where $\tau_i = (\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, \mathbf{o}_{t+1}, r_t, d) \forall t = 1 \dots H-1$. For generalizability and real-world capabilities, we train an LfD method using $\mathcal{D}_{Student}$. At this stage, we use the learned dynamics model, trading faster simulation speed for lower accuracy in the learned model. This is also partially responsible for why $\mathcal{D}_{Student}$ contains suboptimal demonstrations.

3.4 Learning from the Optimized Dataset

Our chosen LfD method is DMfD [41], an off-policy LfD method that learns in state and image spaces. As part of tuning, we employ 100 demonstrations, about two orders of magnitude fewer than the 8000 recommended by the original work. To prevent the policy from overfitting to suboptimal demonstrations in $\mathcal{D}_{Student}$, we disable demonstration-state matching, *i.e.*, resetting the agent to demonstration states and applying imitation reward (see Appendix A.2.5). These were originally proposed [42] as reference state initialization (RSI). These modifications are essential for our LfD implementation, where the demonstrations do not come from an expert.

We use the simulator instead of the learned dynamics model \mathcal{T}_ψ at this stage. This is not because it is computationally infeasible to use the learned model, but because accuracy is important in the final reactive policy. From DMfD, the policy π is parameterized by parameters θ , and learns from data collected in a replay buffer \mathcal{B} . The policy loss contains an advantage-weighted loss \mathcal{L}_A where actions are weighted by the advantage function $A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s})$ and temperature parameter λ . It also contains an entropy component \mathcal{L}_E to promote exploration during data collection. The final policy loss \mathcal{L}_π is a combination of these terms (Eq. 2).

$$\mathcal{L}_A = \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{o} \sim \mathcal{B}} \left[\log \pi_\theta(\mathbf{a} | \mathbf{o}) \exp \left(\frac{1}{\lambda} A^\pi(\mathbf{s}, \mathbf{a}) \right) \right] \quad \mathcal{L}_E = \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{o} \sim \mathcal{B}} [\alpha \log \pi_\theta(\mathbf{a} | \mathbf{o}) - Q(\mathbf{s}, \mathbf{a})]$$

$$\mathcal{L}_\pi = (1 - w_E) \mathcal{L}_A + w_E \mathcal{L}_E, \quad 0 \leq w_E \leq 1 \quad (2)$$

where w_E is a tuneable hyper-parameter. The resulting policy is denoted as π_θ . We pre-populate buffer \mathcal{B} with $\mathcal{D}_{Student}$. Using LfD, we extend from state inputs to image observations, and generalize from v_d to any variation sampled from \mathcal{V} .

4 Experiments

Our experiments are designed to answer the following questions: (1) How does **MAIL** compare to state-of-the-art (SOTA) methods in solving tasks? (Sec. 4.2) (2) How well can **MAIL** solve tasks in the real world? (Fig. 4.2) (3) Can **MAIL** generalize to different n -to- m end-effector transfers? (Sec. 4.3) (4) How do different components of **MAIL** affect performance? (Sec. 4.4)

4.1 Tasks

We experiment with cloth manipulation tasks that are easy for humans to demonstrate but difficult to perform on a robot. We also discuss a simpler rearrangement task with rigid bodies to illustrate generalizability. The tasks are shown in Appendix Fig. 6. We choose a pick-and-place action space, as is common for cloth manipulation [43, 6, 44, 45]. Our action space is 6D (pick and place pose). The end-effectors are pickers in simulation, and a two-finger parallel jaw gripper on the real robot.

CLOTH FOLD: Fold a square cloth in half, along a specified line. **DRY CLOTH:** Pick up a square cloth from the ground and hang it on a plank to dry, variant of [46]. **THREE BOXES:** A simple 2D environment where three boxes of different sizes are randomly placed and need to be moved to designated goal locations. This task is used to illustrate the generalizability of **MAIL** with various n -to- m end-effector transfers, and is not used in the SOTA comparisons. For details on metrics and task variants, see Appendix A.1.

We use particle positions as the state for training dynamics models and trajectory optimization. We record pre-programmed demonstrations for the teacher dataset for each task. For non-LfD and LfD RL training, we use a 32x32 RGB image as the visual observation. The instantaneous reward function, used in learning the policy, is the task performance metric at a given state. Further details on architecture and training are in the supplementary material. In all experiments, we compare each method’s normalized performance, measured at the end of the task given by $\hat{p}(t) = \frac{p(s_t) - p(s_0)}{p_{opt} - p(s_0)}$, where p is the performance metric of state s_t at time t , and p_{opt} is the best performance achievable by the task. We use $\hat{p}(H)$ at the end of the episode ($t = H$).

4.2 SOTA comparisons

Many LfD baselines (Sec. 2) are not directly applicable (they do not handle large difference in action space due to different morphologies). We compare **MAIL** with those LfD baselines that produce a policy with image observations, given demonstrations without actions.

1. SAC-CURL [47]: An image-based RL algorithm that uses contrastive learning and SAC [5] as the underlying RL algorithm. It does not require demonstrations.
2. SAC-DrQ [48]: An image-based RL algorithm that uses a regularized Q-function, data augmentations, and SAC as the underlying RL algorithm. It does not require demonstrations.
3. GNS [49]: A SOTA method that represents cloth as a graph and predicts dynamics using a graph neural network (GNN). It does not require demonstrations but learns dynamics on a random action dataset with particle positions. We run this learned model with a planner [43], provided with full state information.
4. SAC-DrQ-IR: A custom variant of SAC [5] that uses DrQ-based [48] image encoding and a state-only imitation reward (IR) to reach the desired state of the object to be manipulated. It does not imitate actions, as they are unavailable.
5. GAIfo [25]: An adversarial imitation learning algorithm that trains a discriminator on state-state pairs (s, s') from both the demonstrator and agent. This is a popular extension of GAIL [13] that learns the same from state-action pairs (s, a) .
6. GPIL [29] A goal-directed LfD method that uses demonstrations and agent interactions to learn a goal proximity function. This function provides a dense reward to train a policy.

Fig. 3 shows the results. In each environment, the first three columns are demonstration-free baselines, and the last four are LfD methods. **MAIL** outperforms all baselines, in some cases by as much as 24%. For the easier CLOTH FOLD task, the SAC-DrQ baseline came within 11% of **MAIL**.

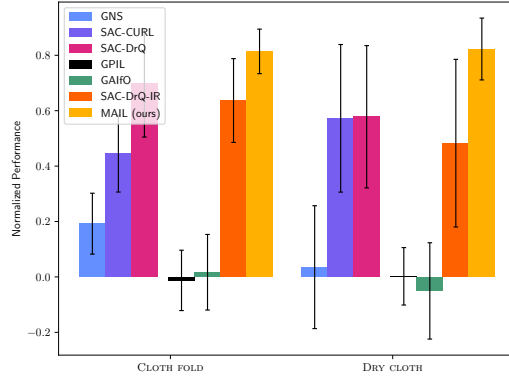


Figure 3: **SOTA performance comparisons.** For each training run, we used the best model in each seed’s training run, and evaluated using 100 roll-outs across 5 seeds, different from the training seed. Bar height denotes the mean, error bars indicate the standard deviation. **MAIL** outperforms all baselines, in some cases by as much as 24%.

267 However, all baselines do not perform well in
 268 the more difficult DRY CLOTH task. RL meth-
 269 ods fail because they have not explored the pa-
 270 rameter space enough without guidance from
 271 demonstrations, and thus converge to a subop-
 272 timal solution. Our custom LfD baseline, SAC-
 273 DrQ-IR, does have reasonable performance, but
 274 the results show that naive imitation alone is not
 275 a good form of guidance to solve it. The other
 276 LfD baselines, GAIfo and GPIL, have poor
 277 performance in both environments. One of the
 278 primary reasons for this is the effect of cross-
 279 morphological demonstrations. They perform
 280 significantly better with student morphology
 281 demonstrations, even if they are suboptimal.
 282 Moreover, environment difficulty also plays an
 283 important part in the final performance. These
 284 and other ablations are described in Sec. 4.4 and
 285 more thoroughly in Appendix Sec. A.2.

286 Surprisingly, the GNS baseline with structured
 287 dynamics does not perform well, even though it
 288 has been used for cloth modeling [50]. We believe that this is because it is designed to learn particle
 289 dynamics via small displacements, but our pick-and-place action space enables large displacements.
 290 Similar to [43], we break down each pick-and-place action into 100 delta actions to work with the
 291 small displacements that GNS is trained on. Thus, planning will accumulate errors from the 100
 292 GNS steps for every action of the planner, which can grow superlinearly due to compounding errors.
 293 This makes it difficult to solve the task. This is especially seen in the DRY CLOTH task (Fig. 3),
 294 where the displacements required to move the entire cloth over the plank are much higher than the
 295 displacements needed for CLOTH FOLD. The rollouts of MAIL on DRY CLOTH show the agent
 296 following the demonstrated guidance - it learned to hang the cloth over the plank. However, it also
 297 displayed an emergent behavior to straighten out the cloth on top of the plank, in an effort to spread
 298 it out to receive higher performance. This was not seen in the two-picker teacher demonstrations.
 299 Demonstrations and rollouts are in the supplementary video file, and on this website.

300 **Real-world results** For DRY CLOTH and CLOTH FOLD tasks, we deploy the learned policies on
 301 a Franka Panda robot (Fig. 5) with a single parallel-jaw gripper. We test the policies with many
 302 different variations of square cloth (size, rotation, translation, thickness, color, and material). For
 303 performance metrics, see Appendix Sec. A.5 The policies achieve $\sim 80\%$ performance, close to the
 304 average performance of our method in simulation, for both tasks.

305 4.3 Generalizability

306 We show, in a simple THREE BOXES task (Fig. 4), how MAIL learns from a demonstrator morphol-
 307 ogy with n end-effectors and deploys to a robot with m end-effectors. Consider a three-picker agent
 308 that solves the task in one pick-place action. It provides the teacher demonstrations, $\mathcal{D}_{Teacher}$. We
 309 transfer them into one-picker or two-picker demonstrations using indirect trajectory optimization
 310 and the learned dynamics model. These will be the optimized datasets that are fed to a downstream
 311 LfD method. In both cases, the LfD method learns to solve the task with a globally optimal 100%
 312 normalized performance. It generalizes from state inputs in the demonstrations to the image inputs
 313 we receive from the environment. Fig. 4 shows the three picker demonstration, a 3-to-2 and 3-to-1
 314 end-effector transfer. We could also do this for the 2-to-1 case, in which a two-picker teacher's
 315 demonstration would take multiple pick-and-place actions to solve the task. Thus, MAIL can solve
 316 a task using n -to- m end-effector transfer with $n > m$, shown here for 3-to-2, 3-to-1, and 2-to-1
 317 cases. It is trivial to perform the transfer for n -to- m with $n \leq m$. One may simply append the
 318 teacher's action space with $m - n$ arms that do no operations. Thus, MAIL is capable of general
 319 n -to- m end-effector transfer.

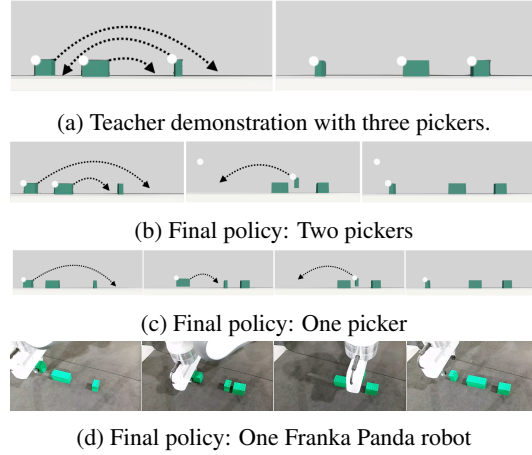


Figure 4: **Sample trajectories of the THREE BOXES task.** A three-picker teacher trajectory to reach the goal state (Fig. 4a). Final policy of the two-picker agent (2 actions to solve the task Fig. 4b). Final policy of the one-picker agent (3 actions to solve the task Fig. 4c). The final policy of the one-picker agent in the real world (Fig. 4d).

4.4 Ablation studies

We use the DRY CLOTH task for all ablations unless specified; it is the most challenging of our tasks. We provide detailed answers to the following questions in Appendix A.2. Appendix Fig. 7 illustrates the ablations corresponding to each part of the overall method. (1) How do different methods perform in creating optimized dataset $\mathcal{D}_{Student}$? (2) What is the best architecture to learn the task dynamics? (3) How good is $\mathcal{D}_{Student}$ compared to the recorded demonstrations? (4) How well does the downstream LfD method handle different kinds of demonstrations? (5) How does the use of expert state matching affect the downstream LfD? (6) How do the baselines perform across related morphologies and environment?

We discovered that the Cross-Entropy Method (CEM) is the most effective optimizer for generating a $\mathcal{D}_{Student}$ from demonstrations. When combined with CEM, the 1D CNN-LSTM architecture produces the best results for trajectory optimization. Our optimized $\mathcal{D}_{Student}$ performs similarly to the pre-programmed \mathcal{D}_{Demo}^{1p} , which has access to full state information of the environment. By utilizing our chosen downstream LfD method, we can successfully complete tasks with a variety of demonstrations and achieve superior performance compared to both $\mathcal{D}_{Student}$ and $\mathcal{D}_{Teacher}$. Expert state matching negatively impacts the performance of DMfD. Lastly, we found that GAIfo trained on our $\mathcal{D}_{Student}$ outperforms GAIfo trained on the $\mathcal{D}_{Teacher}$, and the difficulty of the environment significantly influences the performance of GAIfo and GPIL.

4.5 Limitations

MAIL requires object states in demonstrations and during simulation training, however full state information is not needed at deployment time. It has been tested on the pick-place action space. While it works for high-frequency actions (Appendix A.2.7), it will likely be difficult to optimize actions to create the student dataset for high-dimensional actions. The state-visitation distribution of demonstration trajectories must overlap with that of the student agent; this overlap must contain the equilibrium states of the demonstration. For example, a one-gripper agent cannot reach a demonstration state where two objects are moving simultaneously, but it *can* reach a state where both objects are stable at their goal locations (equilibrium). **MAIL** cannot work when the student robot is unable to reach the goal or intermediate states in the demonstration. For example, in trying to open a flimsy bag with two handles, both end-effectors may simultaneously be needed to keep the bag open. **MAIL** builds a separate policy for each student robot morphology. Subsequent work could learn a single policy conditioned on the desired morphology - another way to think about a base model for generalized LfD.

5 Conclusion

We presented **MAIL**, a framework that enables LfD across morphologies. Our framework enables policy learning for a robot with m end-effectors from teachers with n end-effectors. This enables teachers to record demonstrations in the setting of their own morphology, and vastly expands the set of demonstrations to learn from. We show an improvement of up to 24% over SOTA baselines and discuss other baselines that are unable to handle a large mismatch between teacher and student. Our experiments are on challenging household cloth manipulation tasks performed by a robot with one end-effector based on bimanual demonstrations. We showed that our policy can be deployed zero-shot on a real Franka Panda robot, and generalizes across cloths of varying size, color, material, thickness, and robustness to cloth rotation and translation. We further showed LfD generalizability to any transfer from n -to- m end-effectors, with multiple rigid objects. We believe that this is an important step towards allowing LfD to train a robot to learn from *any* robot demonstrations, regardless of robot morphology, expert knowledge, or the medium of demonstration.

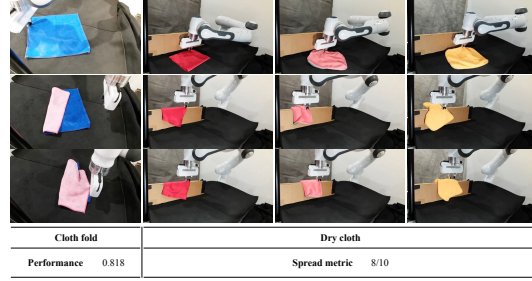


Figure 5: **Real-world results for CLOTH FOLD and DRY CLOTH.** Statistics over 10 rollouts.

References

- [1] K. Pertsch, Y. Lee, Y. Wu, and J. J. Lim. Demonstration-guided reinforcement learning with learned skills. *5th Conference on Robot Learning*, 2021.
- [2] I.-C. A. Liu, S. Uppal, G. S. Sukhatme, J. J. Lim, P. Englert, and Y. Lee. Distilling motion planner augmented policies into visual control policies for robot manipulation. In A. Faust, D. Hsu, and G. Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pages 641–650. PMLR, 08–11 Nov 2022. URL <https://proceedings.mlr.press/v164/liu22b.html>.
- [3] I. Radosavovic, X. Wang, L. Pinto, and J. Malik. State-only imitation learning for dexterous manipulation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7865–7871. IEEE, 2021.
- [4] Y. Yang, Y. Li, C. Fermuller, and Y. Aloimonos. Robot learning manipulation action plans by “watching” unconstrained videos from the world wide web. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1), Mar. 2015. doi:10.1609/aaai.v29i1.9671. URL <https://ojs.aaai.org/index.php/AAAI/article/view/9671>.
- [5] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- [6] J. Hietala, D. Blanco–Mulero, G. Alcan, and V. Kyrki. Learning visual feedback control for dynamic cloth folding. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1455–1462, 2022. doi:10.1109/IROS47612.2022.9981376.
- [7] S. Jin, D. Romeres, A. Raganathan, D. K. Jha, and M. Tomizuka. Trajectory optimization for manipulation of deformable objects: Assembly of belt drive units. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10002–10008. IEEE, 2021.
- [8] J. M. Bern, P. Banzet, R. Poranne, and S. Coros. Trajectory optimization for cable-driven soft robot locomotion. In *Robotics: Science and Systems*, volume 1, 2019.
- [9] D. Pathak, P. Mahmoudieh, G. Luo, P. Agrawal, D. Chen, Y. Shentu, E. Shelhamer, J. Malik, A. A. Efros, and T. Darrell. Zero-shot visual imitation. In *ICLR*, 2018.
- [10] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine. One-shot visual imitation learning via meta-learning. In S. Levine, V. Vanhoucke, and K. Goldberg, editors, *Proceedings of the 1st Annual Conference on Robot Learning*, volume 78 of *Proceedings of Machine Learning Research*, pages 357–368. PMLR, 13–15 Nov 2017. URL <https://proceedings.mlr.press/v78/finnl7a.html>.
- [11] Y. Duan, M. Andrychowicz, B. Stadie, J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba. One-shot imitation learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 1087–1098, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- [12] M. Laskey, J. Lee, R. Fox, A. D. Dragan, and K. Goldberg. DART: noise injection for robust imitation learning. In *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, volume 78 of *Proceedings of Machine Learning Research*, pages 143–156. PMLR, 2017. URL <http://proceedings.mlr.press/v78/laskey17a.html>.
- [13] J. Ho and S. Ermon. Generative adversarial imitation learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/cc7e2b878868cbae992d1fb743995d8f-Paper.pdf>.

- [14] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- [15] M. Vecerík, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. M. O. Heess, T. Rothörl, T. Lampe, and M. A. Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *ArXiv*, abs/1707.08817, 2017.
- [16] A. Mandlekar, J. Booher, M. Spero, A. Tung, A. Gupta, Y. Zhu, A. Garg, S. Savarese, and L. Fei-Fei. Scaling robot supervision to hundreds of hours with roboturk: Robotic manipulation dataset through human reasoning and dexterity. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1048–1055. IEEE, 2019.
- [17] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv e-prints*, pages arXiv–2005, 2020.
- [18] S. Lange, T. Gabel, and M. Riedmiller. Batch reinforcement learning. In *Reinforcement learning*, pages 45–73. Springer, 2012.
- [19] F. Fuchs, Y. Song, E. Kaufmann, D. Scaramuzza, and P. Dürri. Super-human performance in gran turismo sport using deep reinforcement learning. *IEEE Robotics and Automation Letters*, 6(3):4257–4264, 2021.
- [20] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- [21] P. Rashidinejad, B. Zhu, C. Ma, J. Jiao, and S. Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, 2021.
- [22] F. Torabi, G. Warnell, and P. Stone. Recent advances in imitation learning from observation. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 6325–6331. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi:10.24963/ijcai.2019/882. URL <https://doi.org/10.24963/ijcai.2019/882>.
- [23] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 4950–4957. International Joint Conferences on Artificial Intelligence Organization, 7 2018. doi:10.24963/ijcai.2018/687. URL <https://doi.org/10.24963/ijcai.2018/687>.
- [24] Y.-T. A. Sun, H.-C. Lin, P.-Y. Wu, and J.-T. Huang. Learning by watching via key-point extraction and imitation learning. *Machines*, 10(11), 2022. ISSN 2075-1702. doi:10.3390/sun22KeypointExtraction. URL <https://www.mdpi.com/2075-1702/10/11/1049>.
- [25] F. Torabi, G. Warnell, and P. Stone. Generative adversarial imitation from observation. *arXiv preprint arXiv:1807.06158*, 2018.
- [26] J. Yang, J. Zhang, C. Settle, A. Rai, R. Antonova, and J. Bohg. Learning periodic tasks from human demonstrations. *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [27] A. J. Ijspeert, J. Nakanishi, and S. Schaal. Learning rhythmic movements by demonstration using nonlinear oscillators. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, volume 1, pages 958–963, 2002. doi:10.1109/IRDS.2002.1041514.

- [28] F. Al-Hafez, D. Tateo, O. Arenz, G. Zhao, and J. Peters. Ls-ig: Implicit reward regularization for inverse reinforcement learning. In *Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL <https://openreview.net/pdf?id=o3Q4m8jg4BR>.
- [29] Y. Lee, A. Szot, S.-H. Sun, and J. J. Lim. Generalizable imitation learning from observation via inferring goal proximity. In *Advances in Neural Information Processing Systems*, 2021.
- [30] P. T. de Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein. A tutorial on the cross-entropy method. *Annals of Operations Research*, 134:19–67, 2005.
- [31] M. Kelly. An introduction to trajectory optimization: How to do your own direct collocation. *SIAM Review*, 59(4):849–904, 2017.
- [32] M. Posa, C. Cantu, and R. Tedrake. A direct method for trajectory optimization of rigid bodies through contact. *The International Journal of Robotics Research*, 33(1):69–81, 2014.
- [33] Z.-Q. Luo, J.-S. Pang, and D. Ralph. *Mathematical programs with equilibrium constraints*. Cambridge University Press, 1996.
- [34] G. Williams, N. Wagener, B. Goldfain, P. Drews, J. M. Rehg, B. Boots, and E. A. Theodorou. Information theoretic mpc for model-based reinforcement learning. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1714–1721. IEEE, 2017.
- [35] A. Venkatraman, R. Capobianco, L. Pinto, M. Hebert, D. Nardi, and J. A. Bagnell. Improved learning of dynamics models for control. In *2016 International Symposium on Experimental Robotics*, pages 703–713. Springer, 2017.
- [36] T. G. Thuruthel, E. Falotico, F. Renda, and C. Laschi. Learning dynamic models for open loop predictive control of soft robotic manipulators. *Bioinspiration & Biomimetics*, 12(6): 066003, oct 2017. doi:10.1088/1748-3190/aa839f. URL <https://dx.doi.org/10.1088/1748-3190/aa839f>.
- [37] A. S. Polydoros and L. Nalpantidis. Survey of model-based reinforcement learning: Applications on robotics. *Journal of Intelligent & Robotic Systems*, 86(2):153–173, 2017.
- [38] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*. PMLR, 2019.
- [39] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel. Daydreamer: World models for physical robot learning. *Conference on Robot Learning*, 2022.
- [40] M. Macklin, M. Müller, and N. Chentanez. Xpbd: position-based simulation of compliant constrained dynamics. In *Proceedings of the 9th International Conference on Motion in Games*, 2016.
- [41] G. Salhotra, I.-C. A. Liu, M. Dominguez-Kuhne, and G. S. Sukhatme. Learning deformable object manipulation from expert demonstrations. *IEEE Robotics and Automation Letters*, 7(4): 8775–8782, 2022. doi:10.1109/LRA.2022.3187843.
- [42] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 2018.
- [43] X. Lin, Y. Wang, Z. Huang, and D. Held. Learning visible connectivity dynamics for cloth smoothing. In *Conference on Robot Learning*, pages 256–266. PMLR, 2022.
- [44] Y. Avigal, L. Berscheid, T. Asfour, T. Kröger, and K. Goldberg. Speedfolding: Learning efficient bimanual folding of garments. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8, 2022. doi:10.1109/IROS47612.2022.9981402.

- 506 [45] X. Lin, Y. Wang, J. Olkin, and D. Held. Softgym: Benchmarking deep reinforcement learning
507 for deformable object manipulation. *Conference on Robot Learning (CoRL)*, 2020.
- 508 [46] J. Matas, S. James, and A. J. Davison. Sim-to-real reinforcement learning for deformable
509 object manipulation. In *Conference on Robot Learning*, pages 734–743. PMLR, 2018.
- 510 [47] M. Laskin, A. Srinivas, and P. Abbeel. CURL: Contrastive unsupervised representations for re-
511 inforcement learning. In H. D. III and A. Singh, editors, *Proceedings of the 37th International
512 Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*,
513 pages 5639–5650. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/
514 v119/laskin20a.html](https://proceedings.mlr.press/v119/laskin20a.html).
- 515 [48] D. Yarats, I. Kostrikov, and R. Fergus. Image augmentation is all you need: Regularizing
516 deep reinforcement learning from pixels. In *9th International Conference on Learning Repre-
517 sentations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL
518 <https://openreview.net/forum?id=GY6-6sTvGaf>.
- 519 [49] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning
520 to simulate complex physics with graph networks. In *International conference on machine
521 learning*, pages 8459–8468. PMLR, 2020.
- 522 [50] Z. Huang, X. Lin, and D. Held. Mesh-based Dynamics with Occlusion Reasoning for Cloth
523 Manipulation. In *Proceedings of Robotics: Science and Systems*, New York City, NY, USA,
524 June 2022. doi:10.15607/RSS.2022.XVIII.011.

A Appendix

A.1 Tasks

Here we give more details about the tasks, including the performance functions, teacher dataset, and sample images. Fig. 6 shows images all of simulation environments used for SOTA comparisons and generalizability, with one end-effector. In each environment, the end-effectors are pickers (white spheres).

1. CLOTH FOLD: Fold a square cloth in half, along a specified line. The performance metric is the distance of the cloth particles left of the folding line, to those on the right of the folding line. A fully folded cloth should have these two halves virtually overlap. Teacher demonstrations are from an agent with two pickers (*i.e.*, $\mathcal{D}_{Teacher} = \mathcal{D}_{Demo}^{2p}$); we solve the task on a student agent with one picker. Task variations are in cloth rotation.
2. DRY CLOTH: Pick up a square cloth from the ground and hang it on a plank to dry, variant of [46]. The performance metric is the number of cloth particles (in simulation) on either side of the plank and above the ground. Teacher demonstrations are from an agent with two pickers (*i.e.*, $\mathcal{D}_{Teacher} = \mathcal{D}_{Demo}^{2p}$); we solve the task on a student agent with one picker. Task variations are in cloth rotations and translations with respect to the plank.
3. THREE BOXES: A simple 2D environment where three boxes of different sizes are randomly placed and need to be moved to designated goal locations. Teacher demonstrations are from an agent with three pickers (*i.e.*, $\mathcal{D}_{Teacher} = \mathcal{D}_{Demo}^{3p}$); we solve the task on student agents with one picker and two pickers. Performance is measured by the distance of each object from its goal location. This task is used to illustrate the generalizability of **MAIL** with various n -to- m end-effector transfers, and is not used in the SOTA comparisons.

A.2 Ablations

A.2.1 Ablate the method for creating optimized dataset $\mathcal{D}_{Student}$

We answer the question: how do different methods perform in creating optimized dataset $\mathcal{D}_{Student}$? We ablate the optimizer used to create $\mathcal{D}_{Student}$ from the demonstrations, labeled ABL1 in Fig. 7, and compare the following methods, given state inputs from $\mathcal{D}_{Teacher}$.

- Random: A trivial random guesser, that serves as a lower benchmark.
- SAC: An RL algorithm that tries to reach the goal states of the demonstrations.
- Covariant Matrix Adaption Evolution Strategy (CMA-ES): An evolutionary strategy that samples optimization parameters from a multi-variate Gaussian, and updates the mean and covariance at each iteration.
- Cross-Entropy Method (CEM, ours): A well-known gradient-free optimizer, where we assume a Gaussian distribution for optimization parameters.

We did not use gradient-based trajectory optimizers since the contact-rich simulation will give rise to discontinuous dynamics and noisy gradients. As shown in Table 1a, SAC is unable to improve upon the random baseline, likely because of the very large state-space of our environment (> 15000 states for > 5000 cloth particles) and error accumulations from the imprecision of learned dynamics model. Trajectory optimizers achieve the highest performance, and we chose CEM as the best optimizer based on the performance of the optimized trajectory.

A.2.2 Ablate the dynamics model

We answer the question: what is the best architecture to learn the task dynamics? We ablate the learned dynamics model \mathcal{T}_ψ , labeled ABL2 in Fig. 7. The environment state is the state from

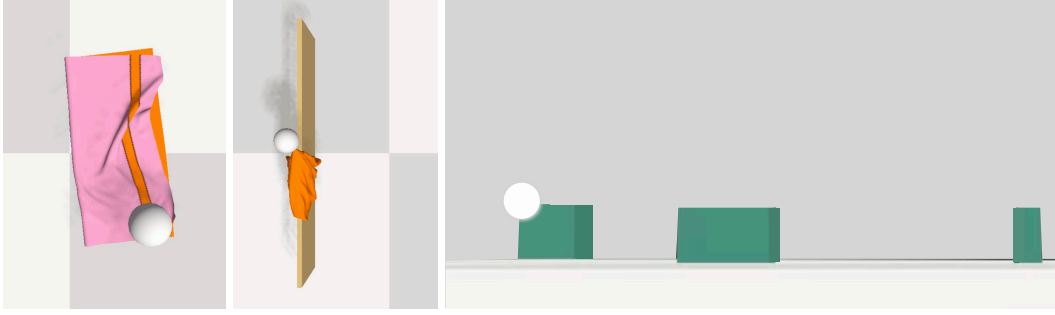


Figure 6: **Environments** used in our experiments, with one end-effector. The end-effectors are pickers (white spheres). In CLOTH FOLD (left) the robot has to fold the cloth (orange and pink) along an edge (inspired by the SoftGym [45] two-picker cloth fold task). In DRY CLOTH (middle) the robot has to hang the cloth (orange and pink) on the drying rack (brown plank). In THREE BOXES (right), the robot has to move three rigid boxes in a 2D environment.

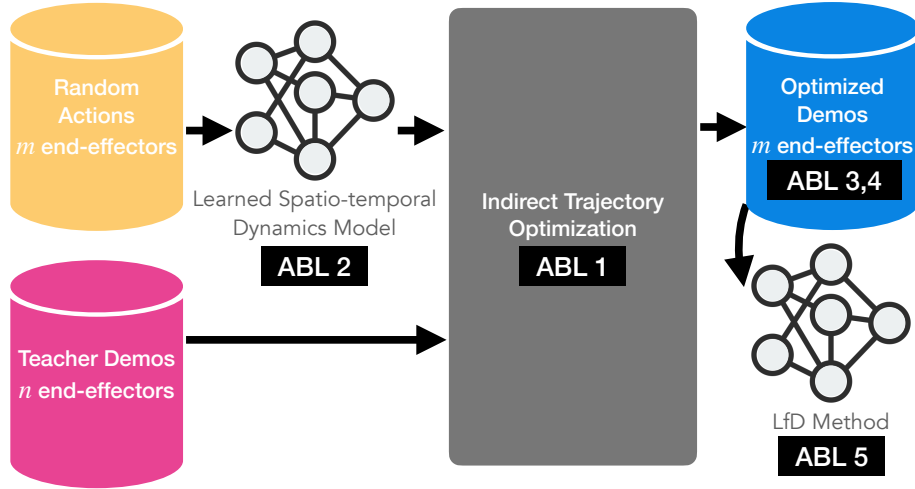


Figure 7: **Ablations** to MAIL components.

570 $\mathcal{D}_{Teacher}$ *i.e.*, positions of cloth particles. This is a structured but large state space since the cloth is
 571 discretized into > 5000 particles.

572 Table 1b shows the performance of trajectories achieved by using the dynamics models. We see that
 573 CNN-LSTM models work better than models that contain only CNNs, graph networks (GNS), or
 574 LSTMs. We hypothesize that this is the case since we need to capture the spatial structure of cloth
 575 and capture a temporal element across the whole trajectory since particle velocity is not captured in
 576 the state. Further, a 1D CNN works better because the cloth state can be simply represented as a 2D
 577 vector ($N \times 3$ which represents the xyz for N particles). This is easier to learn with than the 3D
 578 state vector fed into 2D CNNs.

579 GNS performs poorly also due to the reasons of error accumulation from large displacements, dis-
 580 cussed in Sec. 4.2. Our learned dynamics model \mathcal{T}_ψ was significantly faster than the simulator.
 581 We tested it on a simple training run of SAC [5], without parallelization. Our learned dynamics
 582 gave 162 fps, about $50x$ faster than the 3.4 fps with the simulator. The accuracy was tolerable for
 583 trajectory optimization, as shown in Fig. 8.

584 A.2.3 Compare performance of optimized dataset \mathcal{D}_{Optim}^{1p}

585 We answer the question: how good is $\mathcal{D}_{Student}$ compared to the recorded demonstrations? This
 586 ablation gauges the performance of the optimized dataset that we used as the student dataset for

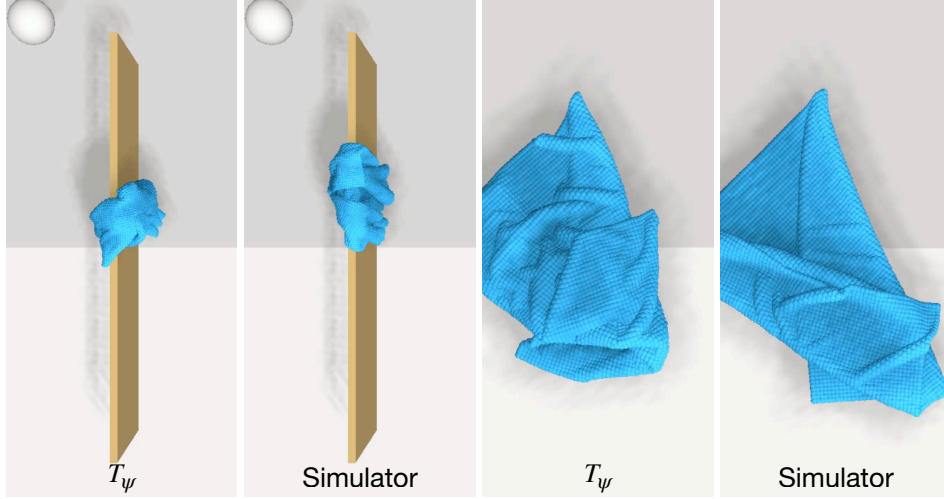


Figure 8: **Predictions of the learned spatio-temporal dynamics model T_ψ and the FleX simulator.** Predictions are made for the same state and action, shown for both cloth tasks. The learned model supports optimization approximately $50\times$ faster than the simulator, albeit at the cost of accuracy.

LfD, $\mathcal{D}_{Student} = \mathcal{D}_{Optim}^{1p}$. We compare this to other relevant datasets to solve the task, as shown in Table 1c. It is labeled ABL3 in Fig. 7. The two-picker demonstrations \mathcal{D}_{Demo}^{2p} are recorded for an agent with two pickers as end-effectors. This is used as the teacher demonstrations in our experiment $\mathcal{D}_{Teacher} = \mathcal{D}_{Demo}^{2p}$. The one-picker demonstrations \mathcal{D}_{Demo}^{1p} are recorded for an agent with one picker as an end-effector. This is to contrast against the optimized demonstrations in the same morphology, \mathcal{D}_{Optim}^{1p} . The random action trajectories are with a one-picker agent, added as a lower performance benchmark. They are the same random trajectories used to train the spatio-temporal dynamics model T_ψ . Naturally, the teacher dataset is the best, as it is trivial to do this task with two pickers. The one-picker dataset has about the same performance as the optimized dataset \mathcal{D}_{Optim}^{1p} , both of which are suboptimal, as it is not trivial to manipulate cloth with one hand. *This is the kind of task we wish to unlock with this work: tasks that are easy to do for teachers in one morphology but difficult to program or record demonstrations for in the student’s morphology.* Note that \mathcal{D}_{Optim}^{1p} has been optimized on the fast but inaccurate learned dynamics model, which is one reason for the reduced performance. This is why the downstream LfD method uses the simulator, as accuracy is very important in the final policy.

A.2.4 Ablate modality of demonstrations

We answer the question: how well does the downstream LfD method handle different kinds of demonstrations? This ablates the composition of the student dataset fed into LfD, and is labeled ABL4 in Fig. 7. We compare the following datasets for $\mathcal{D}_{Student}$, using the notation for datasets explained in Sec. 3.1:

- Demonstrations in one-picker morphology, \mathcal{D}_{Demo}^{1p} : These are non-trivial to create and are thus not as performant, discussed above. Creating these is increasingly difficult as the task becomes more challenging.
- Optimized demos, \mathcal{D}_{Optim}^{1p} : This is optimized from the two-picker teacher demonstrations ($\mathcal{D}_{Teacher} = \mathcal{D}_{Demo}^{2p}$), which are easy to collect as the task is trivial with two pickers.
- 50% \mathcal{D}_{Demo}^{1p} and 50% \mathcal{D}_{Optim}^{1p} : A mix of trajectories from the two cases above. This is an example of handling multiple demonstrators with different morphologies.

Method	25 th %	$\mu \pm \sigma$	median	75 th %
Random	0.000	0.003 \pm 0.088	0.000	0.000
SAC	0.000	0.000 \pm 0.006	0.000	0.000
CMA-ES	0.104	0.270 \pm 0.258	0.286	0.489
CEM	0.351	0.502 \pm 0.242	0.501	0.702

(a) Ablation on the method chosen for creating demonstrations.

Method	25 th %	$\mu \pm \sigma$	median	75 th %
GNS	-0.182	0.002 \pm 0.223	-0.042	0.149
2D CNN, LSTM	0.157	0.376 \pm 0.305	0.382	0.602
No CNN, LSTM	0.327	0.465 \pm 0.213	0.463	0.595
1D CNN, No LSTM	0.202	0.407 \pm 0.237	0.387	0.587
1D CNN, LSTM (ours)	0.351	0.502 \pm 0.242	0.501	0.702

(b) Ablation on the dynamics network architecture.

Dataset	25 th %	$\mu \pm \sigma$	median	75 th %
\mathcal{D}_{Random}	0.000	0.003 \pm 0.088	0.000	0.000
\mathcal{D}_{Demo}^{1p}	0.344	0.484 \pm 0.169	0.446	0.641
\mathcal{D}_{Demo}^{2p}	0.696	0.744 \pm 0.068	0.724	0.785
\mathcal{D}_{Optim}^{1p}	0.351	0.502 \pm 0.242	0.501	0.702

(c) Compare the performance of the optimized dataset.

Table 1: Ablation results for MAIL

Fig. 9 illustrates that all three variants achieve similar final performance. This demonstrates that the downstream LfD method is capable of solving the task with a variety of suboptimal demonstrations. This could be from one dataset of demonstrations, or even a combination of datasets obtained from a heterogeneous set of teachers.

An interesting observation here is that by comparing Fig. 9 and Table 1c, we see that the final policy is better than the suboptimal demonstrations by a considerable margin, and also slightly improves upon the performance of the teacher demonstrations. This improvement comes from the LfD method’s ability to effectively utilize demonstrations and generalize across task variations. This result, combined with the ablation that we need demonstrations in Sec. 4.2, shows that our downstream LfD method is well adapted to work with suboptimal demonstrations to solve a task.

A.2.5 Ablate Reference State Initialization in DMfD

We answer the question: how does the use of demonstration state matching affect the downstream LfD? An improvement we made over the original DMfD algorithm is to disable matching with expert states, known as RSI-IR, first proposed in [42]. We justify this improvement in this ablation, labeled ABL5 in Fig. 7.

As shown in Fig. 10, removing RSI and IR has a net positive effect throughout training, and around 10% on the final policy performance. This means that matching expert states exactly via imitation reward does not help, even during the initial stages of training when the policy is randomly initialized. We believe this is because RSI helps when there are hard-to-reach intermediate states that the policy cannot reach during the initial stages of training. This is true for dynamic or long-horizon tasks, such as karate chops and roundhouse kicks. However, our tasks are quasi-static, and also have a short horizon of 3 for the cloth tasks. In other words, removing this technique allows the policy to

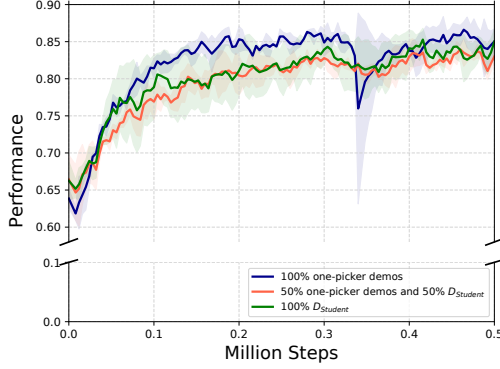


Figure 9: **Ablation on the modality of demonstrations on LfD performance.** Similar performance shows that MAIL can learn from a wide variety of demonstrations, or even a mixture of them, without loss in performance. See Sec. A.2.4.

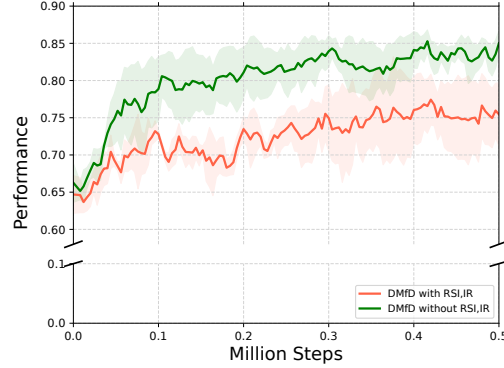


Figure 10: **Ablation on the effect of reference state initialization (RSI) and imitation reward (IR) on LfD performance.** RSI is not helpful here because our tasks are not as dynamic or long horizon as DeepMimic [42]. See Sec. A.2.5.

freely explore the state space while the demonstrations can still guide the RL policy learning via the advantage-weighted loss from DMfD.

A.2.6 Ablate the effect of cross-morphology on SOTA

We answer the question: how do established LfD baselines perform across morphologies? We studied the effect of why baselines such as GAIfO and GPIL performed so poorly on our tasks. In our experiments, we noticed a number of factors (such as variations in the task, diversity of demonstrations, etc.). This ablation studies the effect of cross-morphology in the demonstrations, where we compare the performance of GAIfO, when provided demonstrations from the teacher dataset $\mathcal{D}_{Teacher}$ and student dataset $\mathcal{D}_{Student}$.

As we can see in Table 2, there is a 36% performance improvement when using the (suboptimal) student dataset. Obviously, since the demonstration actions are not available to learn from, the primary difference that the agent sees during training is the richness of demonstration states. Thus, improvement is because of the demonstration states seen in the student dataset. Since the student morphology has only one picker, any demonstration for the task (DryCloth) includes multiple intermediate states of the cloth in various conditions of being partially hung for drying. By contrast, the teacher requires fewer pick-place steps to complete the task, and thus there are fewer intermediate states in the demonstrations.

A.2.7 Ablate the effect of environment difficulty on LfD baselines

We answer the question: how do established LfD baselines perform across environments? Given the subpar performance of the LfD baselines GAIfO and GPIL on our SOTA environments, we ablated the effect of environment difficulty. We took the easy cloth environment (CLOTH FOLD) and used an easier variant of it, CLOTH FOLD DIAGONAL PINNED [41]. In this variant, the agent has to perform an easier fold, but one corner of the cloth is pinned to prevent sliding. Moreover, the desired fold is across the diagonal of the cloth, which can be done by manipulating only one corner of the cloth. We used the state-based observations, and the action space is the small-displacement action space, where the agent outputs incremental picker displacements instead of pick-and-place locations. This action space is similar to those seen in the experiments of GNS, GAIfO and GPIL, where they worked with rigid objects in simulation. This is an easy version of our CLOTH FOLD environment. We can see in Table 3 that the same baselines are able to perform significantly better in this environment. Hence, we believe manipulating with long-horizon pick-place actions, with an

image observation, makes it challenging for the baselines to work in challenging cloth environments described in Sec. 4.1.

Method	25 th %	$\mu \pm \sigma$	median	75 th %
D _{Teacher}	-0.198	-0.055±0.183	-0.043	0.078
D _{Student}	0.199	0.363±0.245	0.409	0.528

Table 2: Ablation of GAIfo on the effect of cross-morphology. We compare the normalized performance, measured at the end of the task.

Method	25 th %	$\mu \pm \sigma$	median	75 th %
GPIL	0.356	0.427±0.162	0.487	0.553
GAIfo	0.115	0.374±0.267	0.471	0.592

Table 3: Measuring performance on the easy cloth task, CLOTH FOLD DIAGONAL PINNED. We compare the normalized performance, measured at the end of the task.

A.3 List of environments we tried for LfD baseline ablations

Two LfD baselines, GAIfo and GPIL, seemed to perform quite poorly, although we expected better performance. In an effort to understand why these fail, we performed a host of studies with different varieties of easier environments, to isolate the properties of the environment that make it the most challenging to succeed. A list of the different task variants we tried are given below. The ones with the most striking difference in performance are described in further detail in Sec. A.2.6 and Sec. A.2.7.

1. Used the easier CLOTH FOLD environment instead of DRY CLOTH.
2. Used state-based environments instead of image-based environments.
3. Reduced the number of variations of the task distribution \mathcal{V} .
4. Used the small-displacement action-space that is used in GNS and GAIfo experiments, instead of the large-displacement pick-place action spaces.
5. Removed the effect of cross-morphology, by providing demonstrations in the students morphology.

A.4 Hyperparameter choices for MAIL

In this section, Table 4 shows the hyperparameters chosen for training the inverse dynamics model \mathcal{T}_ψ . Table 5 shows the details of CEM hyperparameter choices. Table 6 shows the hyperparameters for our chosen LfD method (DMfD).

Parameter	Description
CNN	4 layers, 32 channels, 3x3 kernel, leaky ReLU activation. stride = 2 for the first layer, stride = 1 for subsequent layers
LSTM	One layer Hidden size = 32
Other Parameters	Learning rate $\alpha = 1e-5$ Batch size = 128

Table 4: Hyper-parameters for training the forward dynamics model.

	Planning Horizon	Number of optimization iterations	Number of env interactions
1	1	2	21,000
2	2	2	15,000
3	2	2	21,000
4	2	2	31,000
5	2	2	34,000
6	2	10	21,000
7	2	1	21,000
8	2	1	15,000
9	2	1	32,000
10	3	2	21,000
11	3	10	21,000
12	4	2	21,000
13	4	10	21,000

Table 5: CEM hyper-parameters tested for tuning the trajectory optimization. We conducted ten rollouts for each parameter set and used the set with the highest average normalized performance on the teacher demonstrations. Population size is determined by the number of environment interactions. The number of elites for each CEM iteration is 10% of population size.

Parameter	Description
State encoding	Fully connected network (FCN) 2 hidden layers of 1024, ReLU activation
Image encoding	32x32 RGB input, with random crops. CNN: 4 layers, 32 channels, stride 1, 3x3 kernel, leaky ReLU activation FCN: 1 layer of 1024 neurons, <i>tanh</i> activation
Actor	Fully connected network 2 hidden layers of 1024, leaky ReLU activation
Critic	Fully connected network 2 hidden layers of 1024, leaky ReLU activation
Other parameters	Discount factor: $\gamma = 0.9$ Entropy loss weight: $w_E = 0.1$ Entropy regularizer coefficient: $\alpha = 0.5$ Batch size = 256 Replay buffer size = 600,000 RSI-IR probability = 0 (disabled)

Table 6: Hyper-parameters used in the LfD method (DMfD).

686 A.5 Performance metrics for real-world cloth experiments

687 In this section, we explain the metrics for measuring performance of the cloth, to explain the
688 sim2real results discussed in Fig. 4.2

689 For CLOTH FOLD task, we measure performance at time t by the number of pixels of the top color
690 $pix_{top,t}$ and bottom color $pix_{bot,t}$ of the flattened cloth, compared to the maximum number of pixels,
691 pix_{max} (Fig. 11).

692 For DRY CLOTH task, it is challenging to measure pixels on the sides and top of the plank. Moreover,
693 we could be double counting pixels if they are visible in both side and top views. Hence, we measure
694 the cloth to determine whether the length of the cloth *on top of* the plank is equal to or greater than
695 the side of the square cloth. We call this the spread metric.

696 The policies achieve $\sim 80\%$ performance, which is about the average performance of our method in
697 simulation, for both tasks. However, since these performance metrics are different in the simulation
698 and real world, we cannot *quantify* the sim2real gap through these numbers.

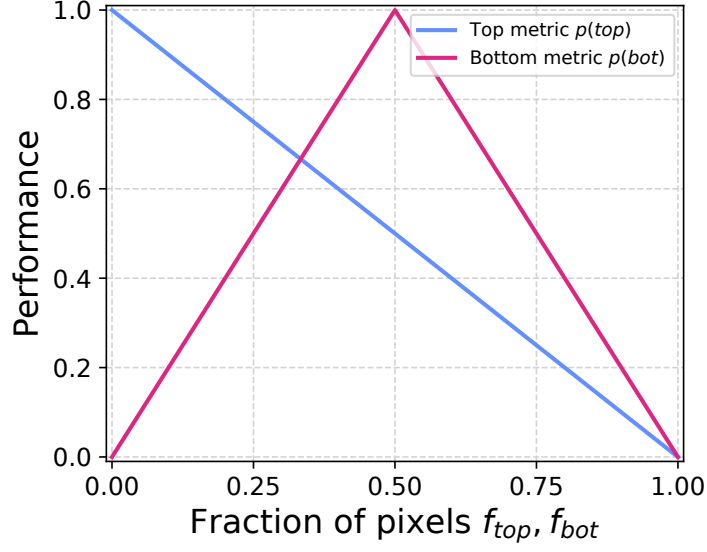


Figure 11: **Performance function for CLOTH FOLD on the real robot.** At time t , we measure the fraction of pixels visible to the maximum number of pixels visible $f_{top} = pix_{top,t}/pix_{max}$ and $f_{bot} = pix_{bot,t}/pix_{max}$. Performance for the top of the cloth should be 1 when it is not visible, $p(top) = 1 - f_{top}$. Performance for the bottom of the cloth should be 1 when it is exactly half-folded on top of the top side, $p(bot) = \min[2(1 - f_{bot}), 2f_{bot}]$. Final performance is an average of both metrics, $p(s_t) = (p(top) + p(bottom))/2$. Note that the cloth is flattened at the start, thus $p_{max} = p_{top,0}$.