


Learning Predictive Representations for Deformable Objects Using Contrastive Estimation

Wilson Yan^{1,†}, Ashwin Vangipuram¹, Pieter Abbeel¹, and Lerrel Pinto^{1,2} 

Abstract: Using visual model-based learning for deformable object manipulation is challenging due to difficulties in learning plannable visual representations along with complex dynamic models. In this work, we propose a new learning framework that jointly optimizes both the visual representation model and the dynamics model using contrastive estimation. Using simulation data collected by randomly perturbing deformable objects on a table, we learn latent dynamics models for these objects in an offline fashion. Then, using the learned models, we use simple model-based planning to solve challenging deformable object manipulation tasks such as spreading ropes and cloths. Experimentally, we show substantial improvements in performance over standard model-based learning techniques across our rope and cloth manipulation suite. Finally, we transfer our visual manipulation policies trained on data purely collected in simulation to a real PR2 robot through domain randomization.

Keywords: Model-based reinforcement learning, contrastive estimation, deformable object manipulation.

1 Introduction

Robotic manipulation of rigid objects has received significant interest over the last few decades, from grasping novel objects in clutter [1, 2] to dexterous in-hand manipulation [3, 4]. However, the objects we interact within our daily lives are not always rigid. From putting on clothes to packing a shopping bag, we constantly need to manipulate objects that deform. Even seemingly rigid objects like metal wires significantly deform during everyday interactions. As a result, there has been a growing interest in algorithms that can tackle deformable object manipulation [5, 6, 7].

Deformable object manipulation presents two key challenges for robots. First, unlike rigid objects, there is no direct representation of the state. Consider the manipulation problem in Figure 1, where the robot needs to straighten a rope from a start configuration to any goal configuration. How does one track the shape of the rope? This lack of a canonical state often limits representations to discrete approximations [8]. Second, the dynamics of deformable objects are complex and non-linear [9]. Due to microscopic interactions within the object, even simple objects can exhibit complex and unpredictable behavior [10], which makes modeling and performing traditional task and motion planning with such deformable objects difficult. One class of techniques that circumvents the challenges in state estimation and dynamics modeling is image-based model-free learning [11]. For instance, Seita et al. [7], Wu et al. [6] use model-free methods in simulation for several difficult cloth manipulation tasks. However, without expert demonstrations, model-free learning is notoriously inefficient [12], and often needs millions of samples to learn from. This challenge is further exacerbated in the multi-task learning framework, where the robot needs to learn to reach multiple goals. Model-based techniques, on the other hand, have shown promise in sample-efficient learning [13, 14]. However, using such model-based learning techniques for deformable objects necessitates tackling the challenges of state representation and dynamics modeling head-on. So how does one learn models given high-dimensional observations and complex underlying dynamics? Some approaches take a direct approach to learning complex dynamics models through pixel-space [15, 16]. Another approach, by Agrawal et al. [17], Nair et al. [18], learns forward dynamics models in conjunction with

^{*1}Department of EECS, University of California, Berkeley.

^{†2}Department of Computer Science, New York University.

^{‡†}Correspondence to wilson1.yan@berkeley.edu

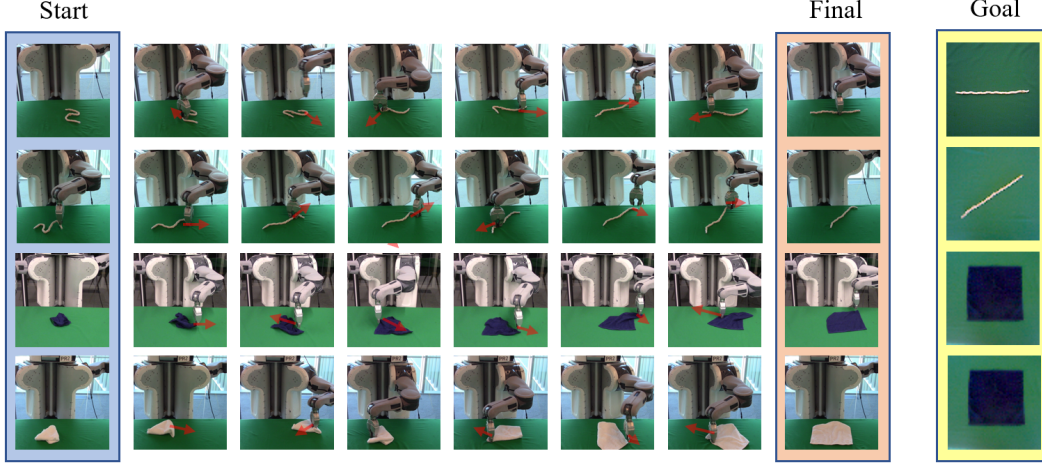


Figure 1: Example trajectories using our contrastive forward model for rope and cloth manipulation. The top two rows show rope manipulation from different start states to different goal states, while the bottom two rows show cloth manipulation using different colored cloths. Note that in the last row, the robot is manipulating a white cloth, but our method is able to still use a blue cloth as the goal image.

inverse dynamic models for manipulating deformable objects. However, during robotic execution, only the inverse model is used. Other model-based approaches such as Wang et al. [19] train Causal InfoGANs [20] to both extract visual representations and forward models, and use the learned forward models for planning. However, these techniques are not robust due to training instabilities associated with GANs [21].

In this paper, we introduce a new visual model-based framework that uses contrastive optimization to jointly learn both the underlying visual latent representations and the dynamics models for deformable objects. We hypothesize that using contrastive methods for model-based learning achieves better generalization and latent space structure due to its inherent information maximization objective. We re-frame the objective introduced in contrastive predictive coding [22] to allow for learning effective model dynamics and latent representations. Once the latent models for representations and dynamics are learned across offline random interactions, we use standard model predictive control (MPC) with one-step predictions to manipulate deformable objects to desired visual goal configurations. Given this controller, we empirically demonstrate substantial improvements over standard model-based learning approaches across multi-goal rope and cloth spreading manipulation tasks.

In summary, we present three key contributions in this paper: (a) We propose a novel contrastive predictive modeling approach to model learning that is compatible with model predictive control. (b) We demonstrate substantial improvements in multi-task deformable object manipulation over other model learning approaches. (c) We show the applicability of our method to real robot rope and cloth manipulation tasks by using sim-to-real transfer without additional real-world training data.

2 Related Work

There has been a substantial amount of prior work in the area of robotic manipulation of deformable objects. A detailed survey of past work can be found in Khalil and Payeur [23], Henrich and Wörn [24]. A standard approach to tackling deformable object manipulation is to use deformable object simulations with planning methods [25]. Past work in this domain has focused on simple linear deformable objects [26, 27], creating better simulations [28], and faster planning [29]. However, the large number of states for deformable objects makes it difficult to plan correctly while being computationally efficient.

Instead of directly planning on the full dynamics, some prior research has focused on planning on simpler approximations, by using local controllers to handle the actual complex dynamics. One approach to using local controllers is model-based servoing [30, 31], where the end-effector is controlled to a goal location instead of explicit planning. However, since the controller is optimized over simple dynamics, it often gets stuck in local minima with more complex dynamics [32]. To

solve this, several works [8] have proposed Jacobian controllers that do not need explicit models, while [33] have proposed learning-based techniques for servoing. We note that our proposed work on learning latent dynamics models is compatible with several of these model-based optimization techniques. Learning good representations remains a difficult challenge in deformable object manipulation. There has been a large amount of prior work on contrastive predictive methods to learn better representations of data. Word2Vec [34] optimizes a contrastive loss to demonstrate semantic and syntactic structure in the learned latent space for words. Oord et al. [22] shows that it is possible to learn high-level representations of images, video, and speech data by employing a large number of negative samples. Tian et al. [35] learns high-level representations by encouraging different views of scenes to be embedded close to one another, and further from others through a similarly framed contrastive loss. Recently, SimCLR [36], another contrastive learning framework, achieved state-of-the-art results in self-supervised learning representations, bridging the gap with supervised learning. There has also been relevant work [37, 38, 39] regarding using contrastive learning to learn better representations in robot manipulation tasks through variants of triplet loss schemes by comparing frames within trajectories. However, none of these works have used contrastive methods for model learning, which is the focus of this work. Kipf et al. [40] uses contrastive learning to learn a world model for object-relation entities. In general, they work with rigid objects, such as shapes, blocks, and Atari objects. However, this would be difficult to apply to learning models for deformable objects which have much more complex dynamics and state representations, especially in cases where objects are interacting with themselves as opposed to other objects.

3 Contrastive Forward Modeling (CFM)

In this section, we describe our proposed framework for learning deformable object manipulation: Contrastive Forward Modeling (CFM). We begin by discussing formalism for predictive modeling and contrastive learning. Following that, we discuss our method for learning contrastive predictive models. See Figure 2 for an overview of our training scheme.

3.1 Dynamic Predictive Models

For our problem setting, we consider a fully observable environment with observations $o \in \mathcal{O}$, actions $a \in \mathcal{A}$, and deterministic transition dynamics $f(o_t, a_t) = o_{t+1}$. We would like to learn a predictive model $\hat{f}(o_t, a_t) \approx o_{t+1}$ to approximate the observation of the next timestep. This can be done by directly learning a visual model through pixel space with regression over observation-action-observation tuples [15]. Once we have successfully learned a predictive model, it is natural to use it for planning to reach different desired goal states, for example, different configurations of a rope or cloth. However, planning directly through pixel space can be difficult, as pixel-value comparisons between images usually do not necessarily correlate well with their true distances. For example, consider an environment with a ball, where the task is to learn a policy that pushes the ball to the center. If the ball is far from the center, then all predicted next actions using a visual forward model would be equidistant from the goal ball-in-center image when comparing pixels values since there would be no image overlap. Therefore, we consider the framework of planning with in a learned latent space by encoding observations. We learn an encoder $g_\theta(o_t) = z_t$ to embed our observations into a latent space, coupled with a predictive model in latent space between z 's, where our learned predictive model is now formulated as $\hat{f}(z_t, a_t) \approx z_{t+1}$. In this work, we propose to learn the latent space using a contrastive learning method.

3.2 Contrastive Models

In our contrastive learning framework, we jointly learn an encoder $g_\theta(o_t) = z_t$ and a forward model $f_\phi(z_t, a_t) \approx z_{t+1}$. We use the InfoNCE contrastive loss described by Oord et al. [22].

$$\mathcal{L} = -\mathbb{E}_{\mathcal{D}} \left[\log \frac{h(\hat{z}_{t+1}, z_{t+1})}{\sum_{i=1}^k h(\hat{z}_{t+1}, \tilde{z}_i)} \right] \quad (1)$$

where h is some similarity function between the computed embeddings from the encoder. The \tilde{z}_i represents negative samples, which are incorrect embeddings of the next state, and we use k such

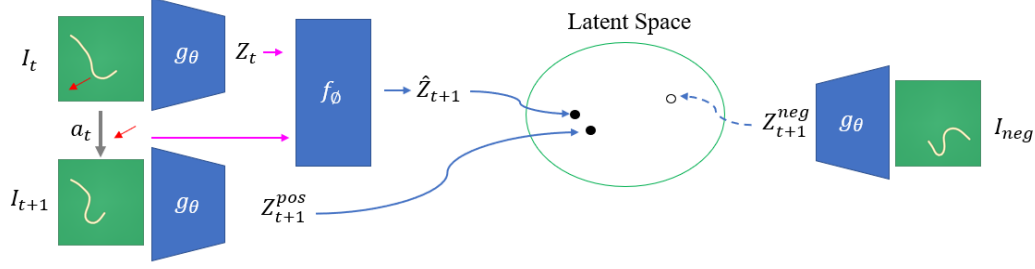


Figure 2: Overview of our contrastive forward model. Training data consists of (image, next image, action) tuples and we learn the encoder and forward model jointly. The contrastive loss objective brings the positive embedding pairs closer together and the negative embeddings further away.

negative samples in our loss. The motivation behind this learning objective lies with maximizing mutual information between the predicted encodings and their respective positive samples. Within the embedding space, this results in the positive sample pairs being aligned together but the negative samples pushed further apart, as seen in Figure 2. Since we are jointly learning a forward model that seeks to minimize $\|f_\phi(z_t, a_t) - z_{t+1}\|^2$, we use the similarity function:

$$h(z_1, z_2) = \exp(-\|z_1 - z_2\|^2) \quad (2)$$

where the norm is a ℓ_2 -norm. After learning the encoder and dynamics model, we plan using a simple version of Model Predictive Control (MPC), where we sample several actions, run them through the forward model from the current z_t , and choose the action a_t that produces \hat{z}_{t+1} closest (in ℓ_2 -distance) to the goal embedding.

4 Experimental Evaluations

In this section, we experimentally evaluate our method in various rope and cloth manipulation settings, both in simulation and in the real world. Our experiments seek to address the following questions:

- Do contrastive learning methods learn better latent spaces and forward models for planning in deformable object manipulation tasks?
- What aspects of our contrastive learning methods contribute the most to performance?
- Can we successfully manipulate deformable objects on a real-world robot?

4.1 Experiment Setup

To simulate deformable objects such as cloth and rope, we used the Deep Mind Control [41] platform with MuJoCo 2.0 [42]. We use an overhead camera that renders $64 \times 64 \times 3$ RGB images as input observations for training our method.

We design the following tasks in simulation:

1. Rope: The rope is represented by 25 geoms in simulation with a four-dimensional action space: the first 2 are the pixel pick point on the rope, and the last 2 are the x, y delta direction to perturb the rope. At the start of each episode, the rope’s state is randomly initialized by applying 120 random actions.

2. Cloth: The cloth is represented by a 9×9 grid of geoms in simulation with a five-dimensional action space: the first 2 are the pixel pick point on the cloth, and the last 3 are the x, y, z delta direction to perturb the cloth. At the start of each episode, the cloth’s state is randomly initialized by applying 50 random actions. In MuJoCo 2.0, the skin of the cloth can be changed by using images taken of a real cloth.

For both rope and cloth environments, we evaluate our method by planning to a desired goal state image and computing the sum of the pairwise geom distances between the achieved and true goal

Table 1: Quantitative comparisons between different model-based learning methods on rope and cloth manipulation tasks. The metric is the sum of pairwise geom distances between the final observation and goal state, where lower distance is more accurate.

	Horizontal	Vertical	Rope 45°	135°	Random	Flat	Cloth Random
Random	4.75	4.93	4.80	4.87	5.73	7.98	10.12
Autoencoder	1.72 ± 0.31	3.24 ± 1.28	2.11 ± 0.51	2.49 ± 0.64	4.308 ± 1.16	3.24 ± 0.29	4.82 ± 0.0
PlaNet	1.81 ± 0.13	3.36 ± 0.78	2.31 ± 0.72	2.38 ± 0.20	3.037 ± 0.24	4.12 ± 0.21	5.06 ± 0.02
Joint Model	2.13 ± 0.66	4.33 ± 0.85	3.88 ± 0.95	4.02 ± 0.85	1.78 ± 0.09	4.24 ± 0.06	4.70 ± 0.03
Visual Model	2.09 ± 0.13	2.65 ± 0.27	2.55 ± 0.34	2.27 ± 0.17	4.77 ± 0.18	2.20 ± 0.05	4.65 ± 0.10
CFM (Ours)	0.58 ± 0.09	3.08 ± 1.19	2.29 ± 1.42	2.24 ± 0.90	1.52 ± 0.10	2.69 ± 0.25	3.97 ± 0.16

	Horizontal	Vertical	Rope (With DR) 45°	135°	Random	Flat	Cloth (With DR) Random
Random	4.75	4.93	4.80	4.87	5.73	7.975	10.12
Autoencoder	3.29 ± 1.08	3.70 ± 1.47	3.19 ± 1.14	3.30 ± 1.14	4.31 ± 1.16	6.26 ± 1.23	7.08 ± 2.22
PlaNet	2.35 ± 0.56	4.06 ± 1.84	3.73 ± 1.66	3.58 ± 1.46	3.04 ± 0.24	8.74 ± 0.55	10.10 ± 1.56
Joint Model	1.01 ± 0.40	2.29 ± 0.10	1.35 ± 0.59	1.82 ± 0.50	1.78 ± 0.09	4.17 ± 0.17	4.64 ± 0.20
Visual Model	3.05 ± 0.45	5.65 ± 0.37	5.37 ± 0.90	5.11 ± 1.04	4.77 ± 0.18	6.64 ± 0.66	6.07 ± 0.52
CFM (Ours)	0.88 ± 0.21	1.20 ± 0.07	0.99 ± 0.07	0.99 ± 0.17	1.38 ± 0.03	3.99 ± 0.15	4.40 ± 0.06

states. We observe that taking an average of 1000 trials suffices to maintain high-confidence evaluation estimates.

Since collecting real-world data on robots is expensive, our method seeks to address this problem by collecting randomly perturbed rope and cloth data in simulation. Using random perturbations allows for a diverse set of deformable objects and interactions for learning the latent space and dynamics model. We collect 4000 trajectories of length 50 for rope (200k samples), and 8000 trajectories of length 50 for cloth (400k samples).

To show the substantial improvements of our model over prior methods, we compare our method against several baselines: a random policy, a visual forward model, an autoencoder trained jointly with a latent dynamics model, PlaNet [43], and a joint dynamics model [17]. In order to ensure that pick points are always on the rope or cloth, we constrain our pick points using a binary segmentation of the observation image computed by RGB thresholding. During planning, all methods use MPC with one-step prediction.

- **Random Policy:** We sample pick actions uniformly over the binary segmentation, and place actions are sampled uniformly random in a unit square centered around the pick location.
- **Visual Forward Model:** We train a forward model similar to Kaiser et al. [15] to perform modeling and planning purely through pixel space.
- **Autoencoder:** We learn a simple latent space model by jointly training a classical autoencoder with a forward dynamics model. The autoencoder learns to minimize the ℓ_2 -distance between reconstructed and actual images [44].
- **PlaNet:** We train PlaNet [43], a stochastic variant of an autoencoder, as another latent space model. PlaNet models a sequential VAE and optimizes a temporal variational lower bound. We do not learn a reward function, as there is no reward in our environments.
- **Joint Dynamics Model:** We jointly learn a forward and inverse model following Agrawal et al. [17].

For consistency across all latent space models, we use a latent size of 8 for both the rope and the cloth environments. For all methods, we sample 100 possible one-step trajectories when performing closed-loop planning. See Fig 3 for example trajectories from each baseline in comparison to our method. In addition, training detail for our method and the baselines can be found in Appendix A.

4.2 Does Using Contrastive Models Improve Performance?

In this section, we compare the results of using our method with those of our baselines, analyzing the advantages and benefits that contrastive models bring over prior methods. Consider a naive baseline where we replace the InfoNCE loss with an MSE loss. This is equivalent to jointly fitting an

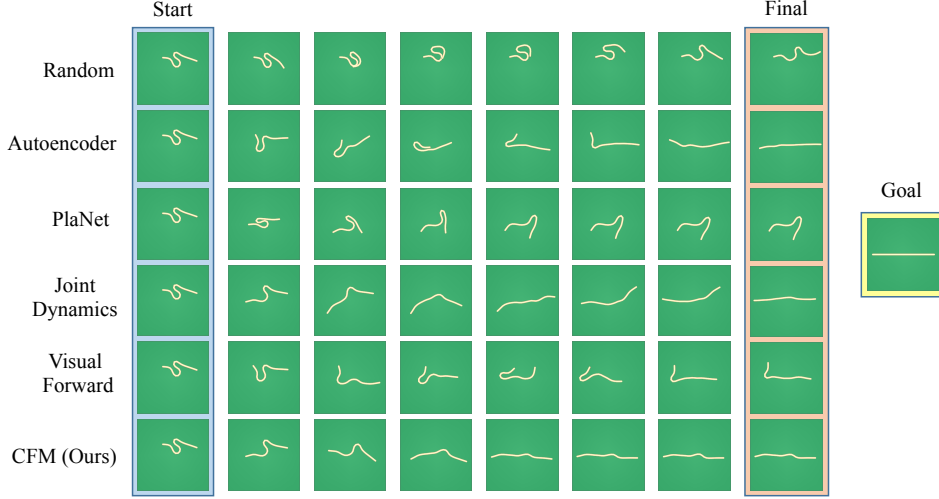


Figure 3: Trajectories for each of the baselines within the simulator, all starting from the same start state and having the same end goal of a horizontal line. Each trajectory was run for 20 actions. Note that our method (CFM) reaches the goal state significantly faster than the baselines.

encoder and dynamics model that minimizes $\|f_\phi(g_\theta(o_t), a_t) - z_{t+1}\|^2$. We can see that the optimal solution is for the encoder to encode all observations to a constant vector to achieve zero loss. To prevent this form of a degenerate solution, we are required to regularize our latent space in some way. Both prior methods and contrastive learning do this in different ways so we analyzed which methods performed better over others. Table 1 shows the quantitative results comparing our method against baselines in different rope and cloth environments, with and without domain randomization for robot transfer. Note that our method does better on all randomly sampled goals with and without domain randomization, indicating stronger generalization in latent spaces for planning. Even in the case where trajectory-wise DR is performed for reconstruction-based baselines, the models show poor generalization to arbitrary goal states. Figure 3 shows example simulator trajectories for each baseline. Each trajectory has the same starting location, same goal image, and was run for 20 actions.

An autoencoder regularizes its latent space by requiring additionally training a decoder to learn to reconstruct o_t from z_t . The model does well in some scenarios, such as a 45° diagonal, but performs poorly when domain randomization is introduced to allow for transfer to a real robot. This is most likely because the autoencoder is optimized to have pixel-level perfect reconstructions, so features such as lighting and color must be encoded in the latent space even when they are not needed for the task. PlaNet behaves similarly to the autoencoder, as it is also a form of a stochastic autoencoder. It performs reasonably competitive with our method but again fails when domain randomization is introduced.

The joint dynamics model regularizes its latent space by jointly learning an inverse model with the forward model. The joint model performed the best across all the baselines when moving to domain randomized data. However, our method still outperforms the joint model for every task.

The visual forward model is the only method that plans in pixel space. It generally performs poorly for tasks with objects with low area coverage, such as the different rope goal orientations, but does better than our method on the cloth flattening task. However, since the forward model operates purely in pixel space, it unsurprisingly suffers from a sharp degradation in performance when introducing domain randomization. As such, it generalizes poorly to the real robot setting.

In addition, we observe that all models generally performed better when given horizontal goals than vertical goals in the rope tasks. This can be explained due to the inherent bias of the environment to initialize rope states to be closer to horizontal than vertical states.

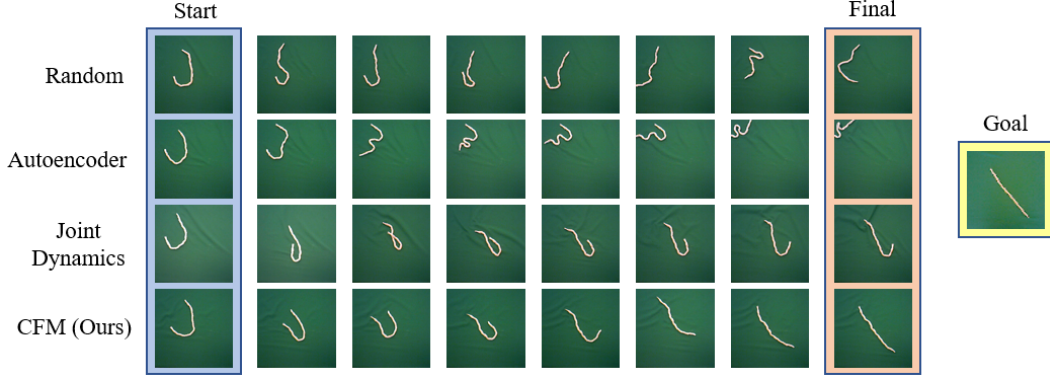


Figure 4: Rope trajectories for each baseline and our method applied on a real robot, all with the same start state, and 135° goal rope orientation.

Table 2: Our results for robot experiments, measuring by the maximum intersection area in pixels between the goal image and observation images averaged over 4 seeds.

Robot Experiments	Rope (Horizontal)	Rope (Vertical)	Rope (45°)	Rope (135°)	Rope (Squiggle)	Cloth (Flat)
Random Policy	6.880	14.727	13.662	4.266	0.049	462.513
Autoencoder	5.526	3.334	3.862	7.499	3.419	603.927
Joint Dynamic Model	17.722	23.636	33.631	21.267	18.311	772.303
CFM (Ours)	32.827	36.387	33.891	38.952	20.711	1001.082

4.3 Ablations on Contrastive Models

In this section, we perform an ablation study on our method, examining the impact of architectural designs on performance. We ablate over two aspects of our method: the forward model architecture, and the contrastive similarity function. For the forward model, our method uses a Multi-Layer Perceptron (MLP) that outputs the parameters of a linear function that is then applied to z_t . For the contrastive similarity function, our method follows Equation 2. The quantitative results, measured as the sum of pairwise geom distances between the final and goal images, appear in Appendix B. First, we compare using our similar function with the original InfoNCE similarity function in Oord et al. [22], the log-bilinear similarity function $h(z_1, z_2) = \exp(z_1^T z_2)$. We achieve the largest boost in performance when switching to our similarity function, as it is more in line with the minimization objective of learning a correct forward model, whereas the log-bilinear model only encourages alignment (as opposed to closeness) of embedding vectors. Lastly, we experiment with a few different forward model architectures: linear, a small MLP, and a small MLP that outputs parameters for a linear transformation. As expected, the biggest drop in performance occurs when learning the simpler linear dynamics model, and a slight drop when using an MLP for both rope and cloth tasks. This demonstrates the need for more complex models for latent forward-dynamics learning.

4.4 Real Robot Experiments

Real Robot Setup We use a PR2 robot to perform our experiments and an overhead camera looking down on the deformable objects to get the RGB image inputs. To ensure the policy learned in the simulator transfers over to the real world, we apply domain randomization by changing the lighting, friction, damping, inertia, and mass of the object during every training step within the simulator. In addition for the cloth environment, we swap between random cloth textures every timestep. We also use a pick and place strategy to mimic the same four-dimensional actions within the simulator. The exact execution process, including data preprocessing, segmentation, and sampling, can be found in Appendix C. Empirically, we found that the execution speed of the robot is bounded by the robot itself as opposed to our method. Executing an action takes around 30-40 seconds whereas planning takes $\ll 1$ second.

Evaluation Metrics We use three baselines along with our contrastive method for real-world evaluation. The first is random actions and the others are the two policies that performed the best with

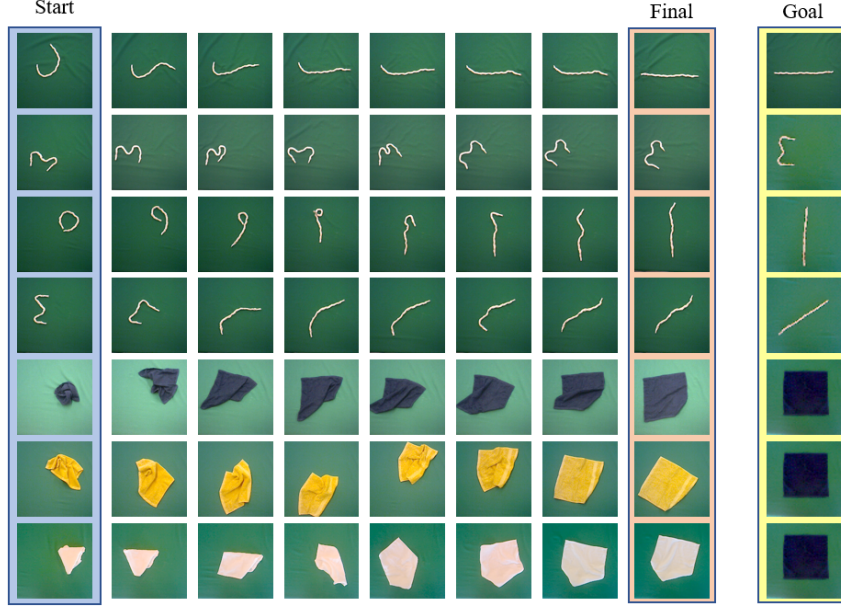


Figure 5: Each row represents one trajectory using our contrastive forward policy. The rope task uses 40 actions between the start and final states, while the cloth task uses 100 actions. Note that for the cloth spreading task the robot is able to use a different colored cloth as a goal.

domain randomization: the autoencoder and the joint dynamics model [17]. For the rope, all the models are evaluated on five goal states: horizontal, vertical, straight line at 45° , straight line at 135° , and a squiggly rope on left. For the cloth, the models are evaluated on one goal state, a flat blue cloth with no rotation. The metric we use is the intersection in pixels between the segmented final image and the segmented goal image. We prefer this instead of intersection over union (IOU) since the objects have the same shape so the union normalization is unnecessary. Additionally, the simpler intersection values provide more insight for comparisons than IOU. The models are run for 40 actions on the rope or 100 actions on the cloth, and the image after each action is stored as an observation. Among all the observations, the one with the highest intersection with the goal is chosen for each method. To account for different seeds, we use 4 starting locations for our contrastive method and 2 starting locations for the baselines, with the scores being averaged across the different start locations. For the cloth, the seed also involves different colors of cloths (blue, gold, white).

The specific evaluation metrics are found in Table 2 which shows that our model performed the best for all the rope and cloth tasks. The joint dynamics model is the second best and got close results to ours on the 45° and squiggle rope tasks. Some example trajectories from our model are seen from a forward view in Figure 1 and from an overhead view in Figure 5. Visual comparisons between our method and baseline methods on the real robot can be seen in Figure 4. We see that our method more accurately plans towards correct goal states compared to the baselines.

5 Conclusion

In this paper, we propose a contrastive learning approach for predictive modeling of deformable objects. We show that contrastive learning learns stronger and more plannable latent representations compared to existing methods. Since our method only requires collecting random data in an environment, it allows for easier transfer to real robots without the need for real-world training.

Acknowledgments

We thank AWS for computing resources. We also gratefully acknowledge the support from Berkeley DeepDrive, NSF, and the ONR Pecase award.

References

- [1] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *ISER*, 2016.
- [2] L. Pinto and A. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. *ICRA*, 2016.
- [3] V. Kumar, E. Todorov, and S. Levine. Optimal control with learned local models: Application to dexterous manipulation. In *ICRA*, 2016.
- [4] M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint*, 2018.
- [5] D. Seita, N. Jamali, M. Laskey, A. K. Tanwani, R. Berenstein, P. Baskaran, S. Iba, J. Canny, and K. Goldberg. Deep transfer learning of pick points on fabric for robot bed-making. *arXiv preprint*, 2018.
- [6] Y. Wu, W. Yan, T. Kurutach, L. Pinto, and P. Abbeel. Learning to manipulate deformable objects without demonstrations. *arXiv preprint*, 2019.
- [7] D. Seita, A. Ganapathi, R. Hoque, M. Hwang, E. Cen, A. K. Tanwani, A. Balakrishna, B. Thananjeyan, J. Ichnowski, N. Jamali, K. Yamane, S. Iba, J. Canny, and K. Goldberg. Deep imitation learning of sequential fabric smoothing policies. *arXiv preprint*, 2019.
- [8] D. Berenson. Manipulation of deformable objects without modeling and simulating deformation. In *IROS*, 2013.
- [9] N. Essahbi, B. C. Bouzgarrou, and G. Gogu. Soft material modeling for robotic manipulation. In *Applied Mechanics and Materials*, 2012.
- [10] P. Pierański, S. Przybył, and A. Stasiak. Tight open knots. *The European Physical Journal E*, 2001.
- [11] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint*, 2018.
- [12] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement learning for continuous control. In *ICML*, 2016.
- [13] 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 *ICRA*, 2017.
- [14] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *ICRA*, 2018.
- [15] L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint*, 2019.
- [16] F. Ebert, C. Finn, S. Dasari, A. Xie, A. Lee, and S. Levine. Visual foresight: Model-based deep reinforcement learning for vision-based robotic control. *arXiv preprint arXiv:1812.00568*, 2018.
- [17] P. Agrawal, A. Nair, P. Abbeel, J. Malik, and S. Levine. Learning to poke by poking: Experiential learning of intuitive physics. *NIPS*, 2016.
- [18] A. Nair, D. Chen, P. Agrawal, P. Isola, P. Abbeel, J. Malik, and S. Levine. Combining self-supervised learning and imitation for vision-based rope manipulation. In *ICRA*, 2017.
- [19] A. Wang, T. Kurutach, K. Liu, P. Abbeel, and A. Tamar. Learning robotic manipulation through visual planning and acting. *RSS*, 2019.
- [20] T. Kurutach, A. Tamar, G. Yang, S. J. Russell, and P. Abbeel. Learning plannable representations with causal infogan. In *NeurIPS*, 2018.

- [21] A. Srivastava, L. Valkov, C. Russell, M. U. Gutmann, and C. Sutton. Veegan: Reducing mode collapse in gans using implicit variational learning. In *NIPS*, 2017.
- [22] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint*, 2018.
- [23] F. F. Khalil and P. Payeur. Dexterous robotic manipulation of deformable objects with multi-sensory feedback-a review. In *Robot Manipulators Trends and Development*. 2010.
- [24] D. Henrich and H. Wörn. *Robot manipulation of deformable objects*. Springer Science & Business Media, 2012.
- [25] P. Jiménez. Survey on model-based manipulation planning of deformable objects. *Robotics and computer-integrated manufacturing*, 2012.
- [26] M. Saha and P. Isto. Manipulation planning for deformable linear objects. *T-RO*, 2007.
- [27] M. Moll and L. E. Kavraki. Path planning for deformable linear objects. *T-RO*, 2006.
- [28] S. Rodriguez, X. Tang, J.-M. Lien, and N. M. Amato. An obstacle-based rapidly-exploring random tree. In *ICRA*, 2006.
- [29] B. Frank, C. Stachniss, N. Abdo, and W. Burgard. Efficient motion planning for manipulation robots in environments with deformable objects. In *IROS*, 2011.
- [30] J. Smolen and A. Patriciu. Deformation planning for robotic soft tissue manipulation. In *Advances in Computer-Human Interactions*, 2009.
- [31] T. Wada, S. Hirai, S. Kawamura, and N. Kamiji. Robust manipulation of deformable objects by a simple pid feedback. In *ICRA*, 2001.
- [32] D. McConachie, M. Ruan, and D. Berenson. Interleaving planning and control for deformable object manipulation. In *ISRR*, 2017.
- [33] B. Jia, Z. Hu, Z. Pan, D. Manocha, and J. Pan. Learning-based feedback controller for deformable object manipulation. *arXiv preprint*, 2018.
- [34] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [35] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. *arXiv preprint*, 2019.
- [36] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint*, 2020.
- [37] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1134–1141. IEEE, 2018.
- [38] D. Dwibedi, J. Tompson, C. Lynch, and P. Sermanet. Learning actionable representations from visual observations. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1577–1584. IEEE, 2018.
- [39] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- [40] T. Kipf, E. van der Pol, and M. Welling. Contrastive learning of structured world models. *arXiv preprint arXiv:1911.12247*, 2019.
- [41] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, et al. Deepmind control suite. *arXiv preprint*, 2018.
- [42] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *IROS*, 2012.

- [43] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent dynamics for planning from pixels. *arXiv preprint*, 2018.
- [44] S. Lange and M. Riedmiller. Deep auto-encoder neural networks in reinforcement learning. In *IJCNN*, 2010.

A Training Details

We used the same encoder architectures for all models. The encoder architecture is a series of 6 2D convolutions of kernel sizes $[3, 4, 3, 4, 4, 4]$, strides $[1, 2, 1, 2, 2, 2]$, and filter sizes $[64, 64, 64, 128, 256, 256]$ respectively. We add Leaky ReLU activation in between each convolutional layer. Finally, the output is flattened and fed into a fully connected layer to produce the latent z . The forward model is a multi-layer perceptron (MLP) with two hidden layers of size 32 which outputs the parameters for a linear transformation on z_t . Specifically for our method (CFM), we use the other batch elements as our negative samples for a total of 127 negative samples per positive pair. More training details for the baseline models can be found in the supplementary materials. Each model was trained on a single NVIDIA TitanX GPU, taking roughly 1 – 2 hours. All of our simulated environments, evaluation metrics, and training code will be publicly released.

For PlaNet, following the original paper, the decoder architecture is a dense layer followed by 4 transposed convolutions with kernel size 4 and stride 2 to upscale to the size of the 64×64 image. The visual forward models follow the same convolutional encoder and decoder architectures as the previous model, where actions are processed by a separate dense layer at each resolution, multiplied channel-wise, and broadcasted spatially. The images and actions were centered and scaled to the range of $[-1, 1]$. We trained all models with batch size 128, learning rate 10^{-3} , and an Adam optimizer for 30 epochs.

B Ablations

Table 3: Ablation experiments on forward model architecture and similarity functions, using the same pair-wise geom MSE evaluation metric as the real robot experiment table in the paper (lower is better).

	Rope	Cloth
Linear	2.30 ± 0.31	4.10 ± 0.03
MLP	1.56 ± 0.10	3.95 ± 0.10
Log-bilinear Similarity	3.25 ± 0.43	4.16 ± 0.15
Ours	1.40 ± 0.02	3.81 ± 0.09

C Robot Execution Details

To compute the actions, we employ a model predictive control (MPC) approach of replanning our action at each time step based on the previous image. We segment the rope/cloth against the background to get the list of valid pick locations of the object. We then generate possible actions by uniformly sampling 100 random deltas in $[-1, 1]$ combined with randomly chosen start locations. We feed these into our forward model along with the encoding of our start image to get the latent encoding for each of the next prospective states. To pick the optimal action, we find the location and delta that minimizes the Euclidean distance from these next states to our goal state and return this action to the robot. The delta from the policy is on the scale of $[-1, 1]$ for both x and y coordinates, and we rescale this to $[-5, 5]$ pixels. On the robot side, we use a learned linear mapping to transform from the image’s pixel values to Cartesian coordinates that the robot uses. To emulate the simulator, the robot’s left arm motion is to go to the start location, go down and close the gripper, move up, move to the new location, move down and open the gripper, where the height of the gripper is hard-coded to some manually tuned value.