

Robustness via Retrying: Closed-Loop Robotic Manipulation with Self-Supervised Learning

Frederik Ebert¹, Sudeep Dasari¹, Alex X. Lee¹, Sergey Levine¹, and Chelsea Finn¹

¹Department of Electrical Engineering and Computer Sciences, UC Berkeley, United States
{febert, sdasari, alexlee_gk, svlevine, cbfinn}@berkeley.edu

Abstract: Prediction is an appealing objective for self-supervised learning of behavioral skills, particularly for autonomous robots. However, effectively utilizing predictive models for control, especially with raw image inputs, poses a number of major challenges. How should the predictions be used? What happens when they are inaccurate? In this paper, we tackle these questions by proposing a method for learning robotic skills from raw image observations, using only autonomously collected experience. We show that even an imperfect model can complete complex tasks if it can continuously retry, but this requires the model to not lose track of the objective (e.g., the object of interest). To enable a robot to continuously retry a task, we devise a self-supervised algorithm for learning image registration, which can keep track of objects of interest for the duration of the trial. We demonstrate that this idea can be combined with a video-prediction based controller to enable complex behaviors to be learned from scratch using only raw visual inputs, including grasping, repositioning objects, and non-prehensile manipulation. Our real-world experiments demonstrate that a model trained with 160 robot hours of autonomously collected, unlabeled data is able to successfully perform complex manipulation tasks with a wide range of objects not seen during training.

1 Introduction

Humans have the ability to learn complex skills, such as manipulating objects, through millions of interactions with their environment during their lifetime. These interactions enable us to acquire a general understanding of the physical world and, notably, do not require significant supervision beyond observation of one’s own actions and their consequences. Hence, self-supervised learning through prediction is an appealing direction of research as it enables intelligent systems to leverage and learn from massive amounts of unlabeled raw data to autonomously acquire general skills. Yet, self-supervised learning systems using predictive models of sensory inputs present a number of challenges: planning needs to account for imperfections in the predictive model and the robot needs a grounded mechanism for evaluating predicted futures. How can we enable systems to plan to perform complex tasks from raw sensory observations, even when the predictions are not always accurate?

Prior work on self-supervised robot learning has enabled robots to learn rudimentary, short-term manipulation skills such as grasping [1, 2], singulation [3], pushing [4, 5], and poking [6]. The question that we are concerned with in this work is: can self-supervised predictive models of raw

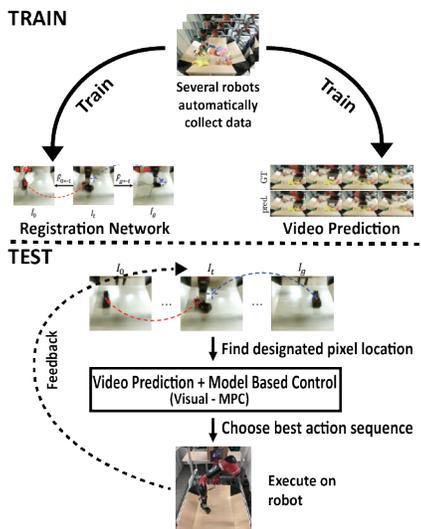


Figure 1: Closed-loop visual MPC: Autonomously collected experience is used to train a video prediction model, as well as an image-to-image registration model, which enables the robot to keep track of its goal.

visual observations be used to perform more complex and realistic tasks, especially tasks that are temporally extended? This might seem challenging due to the difficulty of long-term prediction of image observations. However, if we can enable a robot to continuously correct its behavior during execution and retry in the event of failure or misprediction, even a simple planning algorithm with short-term predictive ability can succeed at the task. This amounts to applying the well-known principle of model-predictive control (MPC), which is typically used with analytical models and low-dimensional state, to control via video prediction. Hence, we call this method visual MPC.

We propose a cost function for vision-based control based on image-to-image registration, which we demonstrate can itself be learned without any supervision using the same exact dataset as the one used to train the predictive model. The key element that enables our method to perform long-horizon tasks is that, using our cost function, the robot can always evaluate the distance to the goal, allowing it to continuously retry, so that even flawed predictions allow for an eventual successful execution.

We demonstrate our method on the task of maneuvering unknown objects in a table-top setting using a robot manipulator. To autonomously learn to perform manipulation skills accurately, tasks need to be specified in way that allows for precision and retrying. We specify a goal by providing an image of the desired configuration along with user-annotated positions of the objects or points of interest¹. This provides a straight-forward and grounded mechanism for a user to specify a goal in the observation space of the robot. Note, however, that these commands are not provided during training, and the model has no prior knowledge of objects or goals. Building on prior methods that use self-supervised predictive models for control [4, 5, 7], we develop a method that can plan actions with a video prediction model to achieve the desired state specified in the goal image.

The main contribution of this work is a method for computing the planning cost for visual MPC based on image-to-image registration, using a learned registration model to put the current observation in correspondence with both the initial and goal images. See Figure 1 for an illustration. This allows closed-loop control, where the robot persistently attempts the task until completion. In contrast to the short-horizon pushing skills demonstrated in prior work [4, 5], we show that our video prediction model can be used to perform longer-horizon manipulations and succeed more consistently. We also demonstrate that this approach can enable a robot to autonomously choose between prehensile and non-prehensile manipulation without any supervision, choosing when to push or pick up objects to relocate them to achieve the goal. Finally, we extend this method to visual MPC with multiple camera videos, showing that task goals can be specified in 3D via prediction from stereo observations.

2 Related Work

While object relocation with both non-prehensile [8, 9] and prehensile [10] manipulation strategies has been explored extensively in prior work, much of this research has focused on hand-designed models and feature representations, which require extensive engineering and make strong assumptions about the environment and task. Video-prediction based manipulation is more general, because it does not require any human intervention at training time, and learns about the physics of object interaction directly from data. A more subtle point is that direct prediction of camera images makes no assumptions about the representation of the world state, such as the number of objects. Many self-supervised robot learning methods have focused on predicting the outcome of a particular event, such as grasp success [1, 2, 3] or crashing [11, 12]. Consequently, these approaches can only recover a single policy that optimizes the probability of the predicted event. Instead, we aim to acquire a model that can be reused for multiple goals and behaviors. To do so, we build upon prior works that learn to predict the sequence of future observations, which can be used to optimize with respect to a variety of goals [4, 5, 7]. Unlike [7], we demonstrate complex object manipulations with previously-unseen objects from RGB video. In contrast to prior self-supervised visual planning works [4, 5], we can perform substantially longer tasks, by using image registration with respect to a goal image.

Goal observations have been previously used for specifying a reward function for robot learning systems [13, 14, 15, 16]. Unlike these methods, we use a learned registration to the goal image for measuring distance to the goal rather than distance in a learned feature space. Distances in unconstrained feature spaces can be arbitrary, while registration inherently measures how pixels should move and can therefore provide a calibrated distance metric with respect to the goal.

¹This also allows the user to specify distractor objects that can be ignored in the goal image.

A related problem is visual servoing, where visual feedback is used to reach a goal configuration [17, 18, 19]. Traditional approaches aim to minimize distances of feature points [20, 21, 22], or pixel intensities [23]. Learning and convolutional networks have also been incorporated into visual servoing [24, 25, 26, 2]. Unlike servoing approaches that use reactive control laws, we use multi-step predictive models to achieve temporally extended goals, while still using continuous visual feedback for retrying at every step. Further, our method performs non-prehensile manipulation, while visual servoing typically assumes fully actuated control, often with a known Jacobian.

Model-predictive control (MPC) [27] has proven successful in a variety of robotic tasks [28, 29, 30, 31, 32]. MPC involves planning with a short horizon and re-planning as the execution progresses, providing the foundation of persistent retrying and the ability to overcome inaccuracies in the predictive model. However, as we later show, maintaining an accurate model of the *cost* used for planning throughout an episode is critical, and has prevented prior work on visual foresight [4, 5] from moving beyond short-term tasks. Our primary contribution is a grounded mechanism for evaluating the planning cost of visual predictions, allowing persistent re-planning with video prediction models.

It is possible to use off-the-shelf trackers [33, 34, 35, 36] to mitigate this issue. However, these trackers usually only have a limited capability to adapt to the domain they are applied to, and can lose track during occlusions. A key advantage of our learned registration approach, inspired by [37], is that we can train our registration model using data collected autonomously by the robot. Since this model is completely self-supervised, it continues to improve as the robot collects more data.

3 Preliminaries

Our visual MPC problem formulation follows the problem statement outlined in prior work [4]. In this setting, an action-conditioned video prediction model g , typically represented by a deep neural network, is used to predict future camera observations $\hat{I}_{1:T} \in \mathbb{R}^{T \times H \times W \times 3}$, conditioned on a sequence of candidate actions $a_{1:T}$, where the prediction horizon is T . This can be written as $\hat{I}_{1:T} = g(a_{1:T}, I_0)$, where I_0 is the frame from the current time-step. An optimization-based planner is used to select the action sequence that results in an outcome that accomplishes a user-specified goal. This type of vision-based control is highly general, since it reasons over raw pixel observations without the requirement for a fixed-size state space, and has been demonstrated to generalize effectively to non-prehensile manipulation of previously unseen objects [4, 5].

Visual MPC assumes the task can be defined in terms of pixel motion. In the initial image I_0 we define n source pixel locations denoted by the coordinates $d_{0,i} \in \mathbb{N}^2$ (for $i \in [0, ..n]$) and the analogous for the goal image I_g denoted by $d_{g,i} \in \mathbb{N}^2$. Given a goal, visual MPC plans for a sequence of actions $a_{1:T}$ to move the pixel at $d_{0,i}$ to $d_{g,i}$. If this pixel lies on top of an object, this corresponds to moving that entire object to a goal position. Note that this problem formulation resembles visual servoing, but it is considerably more complex, since moving the object at d_0 might require complex non-prehensile or prehensile manipulation and long-horizon reasoning. The planning problem is formulated as the minimization of a cost function c , which in accordance with prior work [5], measures the distance between the predicted pixel positions \hat{d}_τ and the goal position d_g for each pixel i :

$$c = \sum_{i=1}^n \lambda_i c_i \quad c_i = \sum_{\tau=1, \dots, T} \mathbb{E}_{\hat{d}_{\tau,i} \sim P_{\tau,i}} \left[\|\hat{d}_{\tau,i} - d_{g,i}\|_2 \right] \quad (1)$$

where $c_i \in \mathbb{R}$ are the costs per source pixel, λ_i are weighting factors discussed in section 4 and $P_{\tau,i}$ is the distribution over predicted pixel positions. The advantage of distance-based cost functions is that they are well-shaped and can be optimized efficiently.

In this paper we use the video prediction model architecture developed by [38], where future images are generated by transforming past images. Starting with a distribution over initial positions of the designated pixel $P_{t_0,i} \in \mathbb{R}^{H \times W}$, $\sum_{H,W} P_{t_0,i} = 1$ at time $t = 0$, the model predicts distributions over its positions $P_{t,i}$ at time $t \in \{1, \dots, T\}$ by exploiting the image transformations used to generate future frames. Planning is performed by sampling candidate actions sequences and optimizing using the cross-entropy method (CEM) [39] to achieve the lowest possible cost c .

To obtain the best results with imperfect models, the action sequence is replanned at each real-world time step² $t \in \{0, \dots, t_{max}\}$ following the framework of model-predictive control (MPC): at each

²We refer to timesteps in the real world as t and to predicted time-steps as τ .

real-world step t , the first action of the best action sequence is executed. At the first real-world time step $t = 0$, the distribution $P_{\tau=0,i}$ is initialized as 1 at the location of the designated pixel and zero elsewhere. In prior work [5, 4], in subsequent steps ($t > 0$), the prediction of the previous step is used to initialize $P_{\tau=0,i}$. However this causes accumulating errors, often preventing the model from solving long-term tasks or responding to situations where the outcome of an action was different than expected. In effect, the model loses track of which object was designated in the initial image.

4 Retrying by Registration

When using a distance-based cost function for visual MPC it is necessary to update the belief of where the target object currently is, so that the agent can “keep retrying” indefinitely until the task is accomplished. Prior work on visual MPC lacked this capability. To solve this issue, we propose a method for registering the current image to both the start and goal image, where the designated pixels are known. In this way we can find the locations of the corresponding pixels in the *current image*, allowing us to compute the distances to the goal. Crucially, the registration method we propose is self-supervised, using the same exact data for training the video prediction model and the registration model. This allows both the predictor and registration model to continuously improve as the robot collects more data.

Before further detailing our learned registration system, we discuss a two simple alternative approaches for obtaining a cost-function for video-prediction based control: One naïve approach could be to use the pixel-wise error between a *goal image* and the *predicted image*. However there are a number of issues with this approach: first, when objects in the image are far from the position in the goal image (e.g., they do not overlap) there is no gradient signal with respect to changes in the actions. Second, due to the blurry predictions from a video prediction model, the pixel-wise difference between the predictions and the goal image can become meaningless.

Another approach is to perform a registration between predicted video frames and the goal image, and use the average length of the warping vectors as a cost function for “closeness” to the goal image. However a major drawback of cost functions based on metrics computed on the complete image is that they naturally emphasize large objects (such as the robot’s arm), while small objects only contribute negligible amounts to the costs. As a result, the planner only tries to match the positions of the large objects (the arm), ignoring smaller objects.

4.1 Test time procedure

We will first describe the registration scheme at test time (see Figure 3(a)). We separately register the current image I_t to the start image I_0 and to the goal image I_g by passing it into the registration network R , implemented as a fully-convolutional neural network. Details about the architecture are given in section 4.3. The registration network produces a flow map $\hat{F}_{0\leftarrow t} \in \mathbb{R}^{H \times W \times 2}$, a vector field with the same size as the image, that describes the relative motion for every pixel between two frames:

$$\hat{F}_{0\leftarrow t} = R(I_t, I_0) \quad \hat{F}_{g\leftarrow t} = R(I_t, I_g) \quad (2)$$

The flow map $\hat{F}_{0\leftarrow t}$ can be used to warp the image of the current time step t to the start image I_0 , and $\hat{F}_{g\leftarrow t}$ can be used to warp from I_t to I_g (see Figure 2 for an illustration):

$$\hat{I}_0 = \hat{F}_{0\leftarrow t} \diamond I_t \quad \hat{I}_g = \hat{F}_{g\leftarrow t} \diamond I_t \quad (3)$$

where \diamond denotes a bilinear interpolation operator that interpolates the pixel value bilinearly with respect to a location (x, y) and its four neighbouring pixels in the image. In essence for a current image $\hat{F}_{0\leftarrow t}$ puts I_t in correspondence with I_0 , and $\hat{F}_{g\leftarrow t}$ puts I_t in correspondence with I_g . As one might expect, warping works better for images that are closer to each other, and sometimes fails

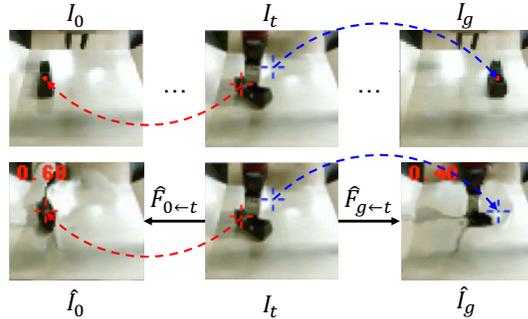


Figure 2: Closed loop control is achieved by registering the current image I_t globally to the first frame I_0 and the goal image I_g . In this example registration to I_0 succeeds while registration to I_g fails since the object in I_g is too far away.

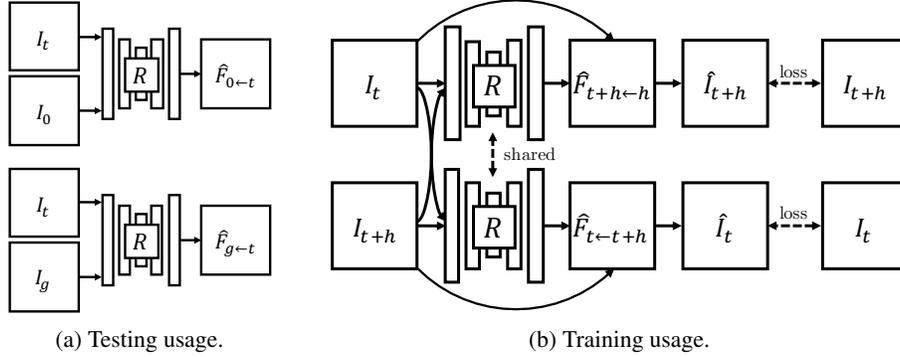


Figure 3: (a) At test time the registration network registers the current image I_t to the start image I_0 (top) and goal image I_g (bottom), inferring the flow-fields $\hat{F}_{0\leftarrow t}$ and $\hat{F}_{g\leftarrow t}$. (b) The registration network is trained by warping images from randomly selected timesteps along a trajectory to each other.

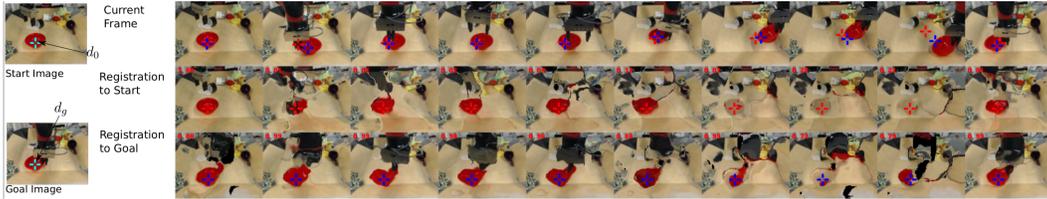


Figure 4: Outputs of registration network. The first row shows the timesteps from left to right of a robot picking and moving a red bowl, the second row shows each image warped to the initial image via registration, and the third row shows the same for the goal image. A successful registration in this visualization would result in images that closely resemble the start- or goal image. In the first row, the locations where the designated pixel of the start image d_0 and the goal image d_g are found are marked with red and blue crosses, respectively. It can be seen that the registration to the start image (red cross) is failing in the second to last time step, while the registration to the goal image (blue cross) succeeds for all time steps. The numbers in red, in the upper left corners indicate the trade off factors λ between the views and are used as weighting factors for the planning cost. (Best viewed in PDF)

when the entities in the image are too far apart. The motivation for registering to both I_0 and I_g is to increase accuracy and robustness. In principle, registering to either I_0 or I_g is sufficient.

While the registration network is trained to perform a global registration between the images, we only evaluate it at the points d_0 and d_g chosen by the user. This results in a cost function that ignores distractors. The flow map produced by the registration network is used to find the pixel locations corresponding to d_0 and d_g in the current frame:

$$\hat{d}_{0,t} = d_0 + \hat{F}_{0\leftarrow t}(d_0) \quad \hat{d}_{g,t} = d_g + \hat{F}_{g\leftarrow t}(d_g) \quad (4)$$

For simplicity, we describe the case with a single designated pixel. In practice, instead of a single flow vector $\hat{F}_{0\leftarrow t}(d_0)$ and $\hat{F}_{g\leftarrow t}(d_g)$, we consider a neighborhood of flow-vectors around d_0 and d_g and take the median in the x and y directions, making the registration more stable. Figure 4 visualizes an example tracking result while the gripper is moving an object.

4.2 Planning Costs

Registration can fail when distances between objects in the images are large. During a trajectory, the registration to the first image typically becomes harder, while the registration to the goal image becomes easier. We propose a mechanism that estimates which image is registered correctly, allowing us to utilize only the successful registration for evaluating the planning cost. This mechanism gives a high weight λ_i to pixel-distance costs c_i associated with a designated pixel $\hat{d}_{i,t}$ that is tracked successfully and a low, ideally zero, weight to a designated pixel where the registration is poor. We propose to use the photometric distance between the true frame and the warped frame evaluated at $d_{0,i}$ and $d_{g,i}$ as an estimate for *local* registration success. A low photometric error indicates that the registration network predicted a flow vector leading to a pixel with a similar color, thus indicating warping success. However this does not necessarily mean that the flow vector points to the correct location. For example, there could be several objects with the same color and the network could

simply point to the wrong object. Letting $I_i(d_i)$ denote the pixel value in image I_i for position d_i , and $\hat{I}_i(d_i)$ denote the corresponding pixel in the image warped by the registration function, we can define the general weighting factors λ_i as

$$\lambda_i = \frac{\|I_i(d_i) - \hat{I}_i(d_i)\|_2^{-1}}{\sum_j^N \|I_j(d_j) - \hat{I}_j(d_j)\|_2^{-1}}. \quad (5)$$

where $\hat{I}_i = \hat{F}_{i \leftarrow t} \diamond I_t$. The MPC cost is computed as the average of the costs c_i weighted by λ_i , where each c_i is the expected distance (see Equation 1) between the registered point $\hat{d}_{i,t}$ and the goal point $d_{g,i}$. Hence, the cost used for planning is $c = \sum_i \lambda_i c_i$. In the case of the single view model and a single designated pixel, the index i iterates over the start and goal image (and $N = 2$).

The proposed weighting scheme can also be used with multiple designated pixels, as used in multi-task settings and multi-view models, which are explained in section 5. The index i then also loops over the views and indices of the designated pixels.

4.3 Training procedure

For registration we use a deep convolutional neural network R which takes in a pair of images and finds correspondences by warping one image to the other. The network is trained on the same data as the video-prediction model, but it does not share parameters with it.³ Our approach is similar to the optic flow method proposed by Meister et al. [37]. However, unlike this prior work, our method computes registrations for frames that might be many time steps apart, and the goal is not to extract optic flow, but rather to determine correspondences between potentially distant images. For training, two images are sampled at random times steps t and $t + h$ along the trajectory and the images are warped to each other in both directions.

$$\hat{I}_t = \hat{F}_{t \leftarrow t+h} \diamond I_{t+h} \qquad \hat{I}_{t+h} = \hat{F}_{t+h \leftarrow t} \diamond I_t \quad (6)$$

The network, which outputs $\hat{F}_{t \leftarrow t+h}$ and $\hat{F}_{t+h \leftarrow t}$ (see Figure 3), is trained to minimize the photometric distance between \hat{I}_t and I_t and \hat{I}_{t+h} and I_{t+h} , in addition to a smoothness regularizer that penalizes abrupt changes in the outputted flow-field. The details of this loss function follow prior work [37]. We found that gradually increasing the temporal distance h between the images during training yielded better final accuracy, as it creates a learning curriculum. The temporal distance is linearly increased from 1 step to 8 steps at 20k SGD steps. In total 60k iterations were taken.

The network R is implemented as a fully convolutional network taking in two images stacked together along the channel dimension. We use three convolutional layers each followed by a bilinear downsampling operation. This is passed into three layers of convolution each followed by a bilinear upsampling operation (all convolutions use stride 1). By using bilinear sampling for increasing or decreasing image sizes we avoid artifacts that are caused by strided convolutions and deconvolutions.

5 Scaling up Visual Model-Predictive Control

Extension to multiple cameras. Prior work has only considered visual MPC with a single camera [4, 5], where objects are manipulated on a plane. To define goals in 3D, we extend visual MPC to include multiple camera views. Since tasks are defined in terms of pixel motion in 2D image space, the combination of multiple 2D tasks with cameras oriented appropriately defines a 3D task. In our experiments, we show that we can define 3D manipulation tasks, such as lifting an object from the table, that would be ambiguous using only a single camera view. The registration method described in the previous section is used separately per view to allow for dynamic retrying and solving temporally extended tasks. The planning costs from each view are combined using weighted averaging where the weights are provided by the registration network (see equation 5).

Combined prehensile and non-prehensile manipulation. In prior work on video-prediction based robotic manipulation [5, 4], the capabilities that emerged out of self-supervised learning were generally restricted to pushing and dragging objects. To enable more complex tasks, we also

³in principle sharing parameters with the video-prediction model might be beneficial, however this is left for future work



Figure 5: Retrying behaviour of our method combining prehensile and non-prehensile manipulation. In the first 4 time instants shown the agent pushes the object. It then loses the object, and decides to grasp it pulling it all the way to the goal. Retrying is enabled by applying the learned registration to both camera views (here we only show the front view).

explore how visual MPC can enable behaviors that include picking and lifting objects for rearrangement. One of the main challenges with this is that random exploration is unlikely to pick up objects a sufficiently large fraction of the time to allow the model to learn grasping skills. To alleviate this challenge, we incorporate a simple “reflex” during data collection, where the gripper automatically closes when the height of the wrist above the table is lower than a small threshold. This reflex is inspired by the palmar reflex observed in infants [40]. With this primitive, about 20% of training trajectories included some sort of grasp on an object. It is worth noting that, other than this reflex, no grasping-specific engineering was applied to the policy allowing a joint pushing and grasping policy to emerge, see figure 5. In our experiments, we evaluate our method using data obtained both with and without the grasping reflex, evaluating both purely non-prehensile and combined prehensile and non-prehensile manipulation.

6 Experiments

Videos and visualizations can be found on our supplementary webpage: <https://sites.google.com/view/robustness-via-retrying>. Our experimental evaluation, conducted using two Sawyer robotic arms, evaluates the ability of our method to learn both prehensile and non-prehensile object relocation tasks entirely through autonomously collected data and self-supervision. In particular, we aim to answer the following questions: (1) How does our MPC approach with self-supervised image registration compare to alternative cost functions, such as off-the-shelf tracking and forward prediction via flow-based models? (2) When the robot can continuously retry a task using goal image registration, how much is the success rate improved for object relocation tasks? (3) Can we learn predictive models that enable both non-prehensile and prehensile object manipulation? In the appendix we also present additional experimental comparisons in a simulated environment.

To train both our prediction and registration models, we collected 20,000 trajectories of pushing motions and 15,000 trajectories with gripper control, where the robot was allowed to randomly move and pick up objects (we use 150 objects for training and 5 objects for testing; see the appendix for details). The data collection process is fully autonomous, requiring human intervention only to replace and change out the objects in front of the robot. The action space consists of Cartesian movements along the x , y , and z axes, and changes in the azimuth orientation of the gripper, while the grasping action is triggered by a primitive as specified in Section 5. For evaluation, we select novel objects that were never seen during training. The evaluation tasks required the robot to move objects in its environment from a starting state to a goal configuration, and performance is evaluated by measuring the distance between the final object position and goal position. In all experiments, the maximum episode length is 50 time steps.

Pushing with retrying. In the first experiment, we aim to evaluate the performance of different visual MPC cost functions, including our proposed self-supervised registration cost. For this experiment, we disable the gripper control, requiring the robot to push objects to the target.

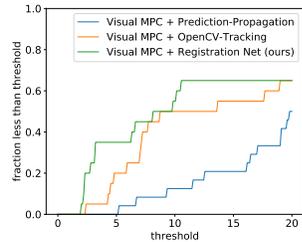


Figure 6: Results for long pushing tasks with 20 objects not seen during training, showing fraction of runs where final distance is lower than threshold. Our method shows a clear gains over OpenCV tracking and predictor propagation.

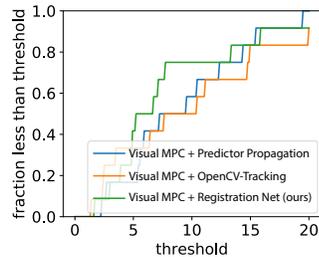


Figure 7: Results for short pushing tasks. Fraction of runs where final distance is lower than threshold.



Figure 8: Applying our method to a pushing task. The goal is to push the bottle to the green point. In the first 3 time instants the object behaves unexpectedly, moving down. The tracking then allows the robot to retry, allowing it to eventually bring the object to the goal.

	Short	Long
Visual MPC + predictor propagation	83%	20%
Visual MPC + OpenCV tracking	83%	45%
Visual MPC + registration network (Ours)	83%	66%

Table 1: Success rate for long-distance pushing benchmark with 20 different object/goal configurations and short-distance benchmark with 15 object/goal configurations. Success is defined as bringing the object closer than 15 pixels to the goal, where the complete image has size 48x64.

We evaluate our method on 20 long-distance and 15 short-distance pushing tasks. For comparisons, we include an ablation where visual-MPC is used with a cost function that computes object locations using the “multiple instance learning tracker” (MIL) [35] in OpenCV. This method represents a strong baseline trained with additional human supervision, which our method does not have access to. We also compare to the visual MPC method proposed by Ebert et al. [5], which does not track the object explicitly, but relies on the flow-based video prediction model to keep track of the designated pixel, which we call “predictor propagation.”

Figure 6 illustrates that, on the long-distance benchmark, where the average distance between the object and its goal position is 30 cm, our approach not only outperforms prior work [5], but also outperforms the hand-designed, supervised object tracker [35]. Using our learned registration, the robot is more frequently able to successfully recover after mispredictions or occlusions (see Figure 8). In contrast, in the short distance benchmark, where the average push distance is 15 cm all methods perform comparably, as shown in Figure 7. These results indicate the importance of closed loop control in long-horizon tasks.

Combined prehensile and non-prehensile manipulation.

In the setting where the grasping reflex is enabled, the robot needs to decide whether to solve a task by grasping or pushing. We again use 20 object relocation tasks to measure the performance of each method, in terms of the final distance between the object and its goal location. We observe that, in the majority of the cases, the robot decides to grasp the object, as can be seen in the supplementary video. Figure 9 shows the results of a benchmark on long-distance relocation tasks, which show that visual MPC combined with our registration method is comparable with the performance of visual MPC with the hand-designed and supervised OpenCV tracker, without requiring any human supervision.

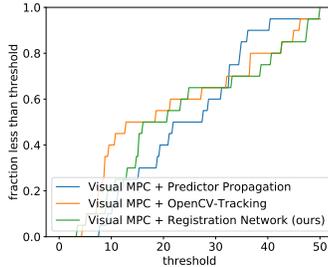


Figure 9: Results with combined prehensile and non-prehensile manipulation.

7 Discussion

We demonstrate that using a cost function derived from learned image-to-image registration substantially improves performance on temporally extended object manipulation tasks by closing the control loop. Our experiments show a large improvement in success rate compared to a prior open-loop visual MPC method [5]. We further show that, by including a simple grasping “reflex” inspired by the palmar reflex in infants, we can efficiently learn both non-prehensile and prehensile object relocation skills in a purely self-supervised way, allowing the robot to plan to pick up and move objects when necessary. Interesting challenges for future work are multi-object relocation tasks, which could be addressed by adding an abstract long-term planner on top of the presented framework.

References

- [1] L. Pinto and A. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *International Conference on Robotics and Automation (ICRA)*, 2016.
- [2] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *International Journal of Robotics Research (IJRR)*, 2016.
- [3] A. Zeng, S. Song, S. Welker, J. Lee, A. Rodriguez, and T. Funkhouser. Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. *arXiv preprint arXiv:1803.09956*, 2018.
- [4] C. Finn and S. Levine. Deep visual foresight for planning robot motion. In *International Conference on Robotics and Automation (ICRA)*, 2017.
- [5] F. Ebert, C. Finn, A. X. Lee, and S. Levine. Self-supervised visual planning with temporal skip connections. *CoRR*, abs/1710.05268, 2017. URL <http://arxiv.org/abs/1710.05268>.
- [6] P. Agrawal, A. V. Nair, P. Abbeel, J. Malik, and S. Levine. Learning to poke by poking: Experiential learning of intuitive physics. In *Advances in Neural Information Processing Systems*, 2016.
- [7] A. Byravan, F. Leeb, F. Meier, and D. Fox. SE3-Pose-Nets: Structured deep dynamics models for visuomotor planning and control. *arXiv preprint arXiv:1710.00489*, 2017.
- [8] T. Hermans, F. Li, J. M. Rehg, and A. F. Bobick. Learning contact locations for pushing and orienting unknown objects. Georgia Institute of Technology, 2013.
- [9] M. Salganicoff, G. Metta, A. Oddera, and G. Sandini. *A vision-based learning method for pushing manipulation*, volume 54. University of Pennsylvania, 1993.
- [10] C. Goldfeder, M. Ciocarlie, J. Peretzman, H. Dang, and P. K. Allen. Data-driven grasping with partial sensor data. In *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, pages 1278–1283. Citeseer, 2009.
- [11] D. Gandhi, L. Pinto, and A. Gupta. Learning to fly by crashing. *arXiv preprint arXiv:1704.05588*, 2017.
- [12] G. Kahn, A. Villaflor, V. Pong, P. Abbeel, and S. Levine. Uncertainty-aware reinforcement learning for collision avoidance. *arXiv preprint arXiv:1702.01182*, 2017.
- [13] M. Jagersand and R. Nelson. Visual space task specification, planning and control. In *International Symposium on Computer Vision*, 1995.
- [14] K. Deguchi and I. Takahashi. Image-based simultaneous control of robot and target object motions by direct-image-interpretation method. In *International Conference on Intelligent Robots and Systems (IROS)*, 1999.
- [15] M. Watter, J. Springenberg, J. Boedecker, and M. Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in neural information processing systems*, 2015.
- [16] C. Finn, X. Y. Tan, Y. Duan, T. Darrell, S. Levine, and P. Abbeel. Deep spatial autoencoders for visuomotor learning. In *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016.
- [17] S. Hutchinson, G. D. Hager, and P. I. Corke. A tutorial on visual servo control. *IEEE transactions on robotics and automation*, 12(5):651–670, 1996.
- [18] D. Kragic and H. I. Christensen. Survey on visual servoing for manipulation. *Computational Vision and Active Perception Laboratory, Fiskartorpsv*, 15, 2002.
- [19] G. N. DeSouza and A. C. Kak. Vision for mobile robot navigation: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 24(2):237–267, 2002.
- [20] J. T. Feddema and O. R. Mitchell. Vision-guided servoing with feature-based trajectory generation (for robots). *IEEE Transactions on Robotics and Automation*, 5(5):691–700, 1989.
- [21] B. Espiau, F. Chaumette, and P. Rives. A new approach to visual servoing in robotics. *IEEE Transactions on Robotics and Automation*, 8(3):313–326, 1992.

- [22] W. J. Wilson, C. C. Williams Hulls, and G. S. Bell. Relative end-effector control using cartesian position based visual servoing. *IEEE Transactions on Robotics and Automation*, 12(5):684–696, 1996.
- [23] G. Caron, E. Marchand, and E. M. Mouaddib. Photometric visual servoing for omnidirectional cameras. *Autonomous Robots*, 35(2-3):177–193, 2013.
- [24] A. Saxena, H. Pandya, G. Kumar, A. Gaud, and K. M. Krishna. Exploring convolutional networks for end-to-end visual servoing. In *International Conference on Robotics and Automation (ICRA)*, 2017.
- [25] Q. Bateux, E. Marchand, J. Leitner, and F. Chaumette. Visual servoing from deep neural networks. In *International Conference on Robotics and Automation (ICRA)*, 2018.
- [26] A. X. Lee, S. Levine, and P. Abbeel. Learning visual servoing with deep features and fitted Q-iteration. *International Conference on Learning Representations (ICLR)*, 2017.
- [27] E. F. Camacho and C. B. Alba. *Model predictive control*. Springer Science & Business Media, 2013.
- [28] D. H. Shim, H. J. Kim, and S. Sastry. Decentralized nonlinear model predictive control of multiple flying robots. In *Decision and control, 2003. Proceedings. 42nd IEEE conference on*, volume 4, pages 3621–3626. IEEE, 2003.
- [29] G. Allibert, E. Courtial, and F. Chaumette. Predictive control for constrained image-based visual servoing. *IEEE Transactions on Robotics*, 26(5):933–939, 2010.
- [30] T. M. Howard, C. J. Green, and A. Kelly. Receding horizon model-predictive control for mobile robot navigation of intricate paths. In *Field and Service Robotics*, pages 69–78. Springer, 2010.
- [31] 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 *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 1714–1721. IEEE, 2017.
- [32] I. Lenz and A. Saxena. Deepmpc: Learning deep latent features for model predictive control. In *In RSS*. Citeseer, 2015.
- [33] B. D. Lucas, T. Kanade, et al. An iterative image registration technique with an application to stereo vision. 1981.
- [34] T. Brox, A. Bruhn, N. Papenberg, and J. Weickert. High accuracy optical flow estimation based on a theory for warping. In *European conference on computer vision*. Springer, 2004.
- [35] B. Babenko, M.-H. Yang, and S. Belongie. Visual tracking with online multiple instance learning. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 983–990. IEEE, 2009.
- [36] X. Mei and H. Ling. Robust visual tracking using l1 minimization. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 1436–1443. IEEE, 2009.
- [37] S. Meister, J. Hur, and S. Roth. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. *arXiv preprint arXiv:1711.07837*, 2017.
- [38] A. X. Lee, R. Zhang, F. Ebert, P. Abbeel, C. Finn, and S. Levine. Stochastic adversarial video prediction. *arXiv preprint arXiv:1804.01523*, 2018.
- [39] R. Y. Rubinstein and D. P. Kroese. *The cross-entropy method: a unified approach to combinatorial optimization, Monte-Carlo simulation and machine learning*. Springer Science & Business Media, 2013.
- [40] D. Sherer. Fetal grasping at 16 weeks’ gestation. *Journal of ultrasound in medicine*, 12(6), 1993.
- [41] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 5026–5033. IEEE, 2012.

Appendix

Simulated Experiments

In order to provide a more controlled comparison, we also set up a realistic simulation environment using MuJoCo [41], which includes a robotic manipulator controlled via Cartesian position control, similar to our real world setup, pushing randomly-generated L-shaped objects with random colors (see details in supplementary materials). We trained the same video prediction model in this environment, and set up 50 evaluation tasks where blocks must be pushed to target locations with maximum episode lengths of 120 steps.

We compare our proposed registration-based method, “predictor propagation,” and ground-truth registration obtained from the simulator, which provides an oracle upper bound on registration performance. Figure 11 shows the results of this simulated evaluation, where the x-axis shows different distance thresholds, and the y-axis shows the fraction of evaluation scenarios where each method pushed the object within that threshold. We can see that, for thresholds around 0.1, our method drastically outperforms predictor propagation (i.e., prior work [5]), and has a relatively modest gap in performance against ground-truth tracking. This indicates that our registration method is highly effective in guiding visual MPC, despite being entirely self-supervised.

In our simulated experiments, we use end-effector position control with the arm illustrated in Figure 10. In this environment, the video prediction model was trained with using 60,000 training trajectories.

Planner Implementation Details

We use the cross-entropy method [39] to optimize the action sequence with respect to the cost function. At the first iteration 400 samples and at later iterations 200 samples are taken and a Gaussian is fitted to the best 5% of the samples. We use 3 CEM iterations. The length of the prediction horizon is 15, to reduce the search space we use an action repeat of 3, so that only a sequence of 5 independent actions needs to be optimized. Since the action space is 4 ($x, y, z, \text{rotation}$) a total of 20 variables are optimized.

Improvements to online optimization procedure

In the visual MPC setting the action sequences found by the optimizer can be very different between execution time steps (not to be confused with prediction time steps). For example at one time step the optimizer might find a pushing action leading towards the goal and in the next time step it determines a grasping action to the optimal to reach the goal. Naive replanning at every time step can result in alternating between a pushing attempt and a grasping attempt indefinitely causing the agent to get stuck and not making any progress towards to goal.

We show that we can resolve this problem by modifying the sampling distribution of the first iteration of CEM so that the optimizer commits to the plan found in the previous time step. In prior work [5] the sampling distribution at first iteration of CEM is chosen to be a Gaussian with diagonal covariance matrix and zero mean. We instead use the best action sequence found in the optimization of the previous time step as the mean. Since this action sequence is optimized for the previous time step we only use the values $a_{1:T}$ and omit the first action, where T is the prediction horizon. To sample actions close to the action sequence from the previous time step we reduce the entries of the diagonal covariance matrix for the first $T - 1$ time steps. It is crucial that the last entry of the covariance matrix at the end of the horizon is not reduced otherwise no exploration could happen for the last time step causing poor performance at later time steps.

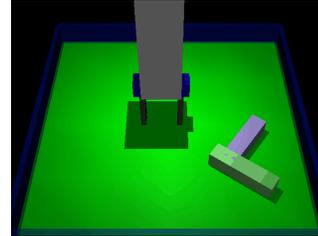


Figure 10: Block pushing simulator

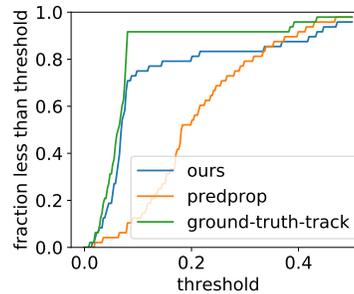


Figure 11: Simulated evaluation. Fraction of trajectories with final object distance lower than threshold (higher is better).