# Improvisation through Physical Understanding: Using Novel Objects as Tools with Visual Foresight

Annie Xie, Frederik Ebert, Sergey Levine, Chelsea Finn
UC Berkeley

*Abstract*—Machine learning techniques have enabled robots to learn narrow, yet complex tasks and also perform broad, yet simple skills with a wide variety of objects. However, learning a model that can both perform complex tasks and generalize to previously unseen objects and goals remains a significant challenge. We study this challenge in the context of "improvisational" tool use: a robot is presented with novel objects and a user-specified goal (e.g., sweep some clutter into the dustpan), and must figure out, using only raw image observations, how to accomplish the goal using the available objects as tools. We approach this problem by training a model with both a visual and physical understanding of multi-object interactions, and develop a sampling-based optimizer that can leverage these interactions to accomplish tasks. We do so by combining diverse demonstration data with self-supervised interaction data, aiming to leverage the interaction data to build generalizable models and the demonstration data to guide the model-based RL planner to solve complex tasks. Our experiments show that our approach can solve a variety of complex tool use tasks from raw pixel inputs, outperforming both imitation learning and self-supervised learning individually. Furthermore, we show that the robot can perceive and use novel objects as tools, including objects that are not conventional tools, while also choosing dynamically to use or not use tools depending on whether or not they are required. Videos of the results are available online[1].

## I. INTRODUCTION

An understanding of physical cause-and-effect relationships is a powerful means for enabling robots to achieve a wide variety of complex goals. This understanding becomes especially useful when performing complex multi-object manipulation tasks, such as those involved in tool use: if a robot could predict how one object might interact with another, it would be able to autonomously construct tool-use behaviors on the fly. While fully-specified analytic and symbolic models of physics can allow fully observable systems to perform such tasks [45], acquiring such models is substantially more challenging when the environment can only be observed through image observations. Learning predictive models of low-level observations, such as camera image pixels, has a number of benefits. Such models can be learned from real world data and deployed in real world settings, as they do not require direct access to the state of the objects in the world. Models from pixels further do not need knowledge of object shapes, surface friction, or other properties, and hence can use large datasets of experience and readily generalize to new objects. Indeed, model-based reinforcement learning with action-conditioned video prediction models, known as *visual*
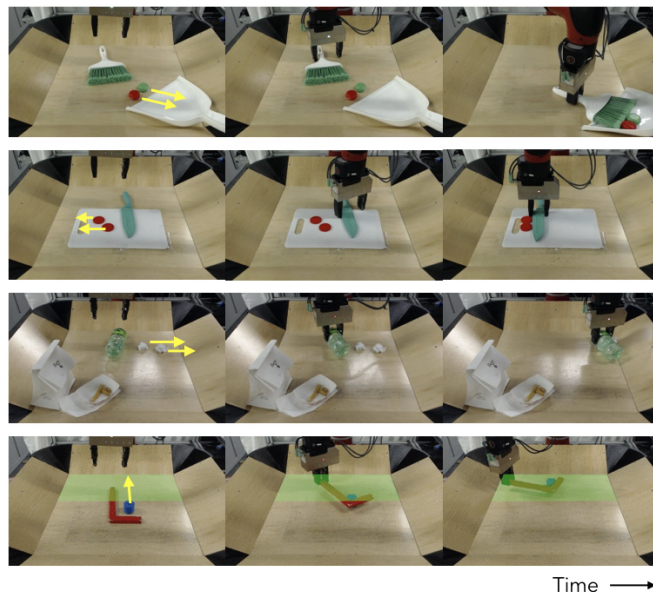
Fig. 1. The robot learns a visual predictive model and uses it to manipulate new objects, that were never seen before, as tools to accomplish tasks specified by a person (as indicated by the yellow arrows). The robot can even utilize objects that are not conventional tools, such as water bottles. In the last example, when the robot is allowed to move only within the shaded green region, it uses the L-shaped hook to pull the blue cylinder towards itself.

*MPC, or visual foresight*, has enabled robots to perform short-horizon tasks involving a range of novel, previously-unseen objects [16, 13, 11].

In this paper, we investigate how models that predict low-level percepts, such as future camera images, can enable a robot to reason about multi-object interactions. In particular, we aim to study "improvisational" tool use, where the robot can use new objects, which might never have been seen before, to interact with other items to perform tasks that are not achievable with only basic single-object manipulation skills (e.g., grasping and pushing with the robot's own end-effector). We focus on tool use in this work because it represents one of the most fundamental multi-object manipulation capabilities, and is a skill that is often associated with greater levels of intelligence in animals [34].

Imitation learning approaches have enabled robots to learn to perform complex tasks [4, 3], including some types of tool use [38], while visual MPC has enabled robots to perform skills with many different novel objects [16, 13]. However, both have their limitations: policies learned through imitation are typically inflexible, as they are constrained to imitate the

demonstrator, while work on visual MPC [16, 13] has so far been limited to simple, short-term skills. We combine ideas from imitation learning and from visual MPC to show that their combination can outperform each approach individually when applied to problems requiring tool use. In particular, we show that we can use demonstrations to solve complex tasks, while retaining the flexibility of visual planning.

The main contribution of this paper is a study of how direct prediction of low-level sensory observations, namely camera images, can enable a robot to carry out improvised multi-object interactions – that is, determine how to use tools in its environment to perform tasks that require tool use. To this end, we combine ideas from imitation learning and prior work on visual MPC [16], incorporating imitation-driven models into both the data collection process and the sampling-based planning procedure. Our method uses video prediction to reason about potential robot actions, constructing plans to manipulate novel objects on the fly, in less than a second. Our experiments with a real-world Sawyer robot indicate that, by leveraging demonstrations and autonomously collected, self-supervised data, the robot can decide to use tools in situations where they are needed and use its arm directly, without the tool, in situations where tools are not needed. Further, by reasoning about object interactions, the robot can find effective tool-use strategies even if it has never seen the tool before, and even in situations where no conventional tools are available. In comparisons, our approach exceeds the performance of both direct imitation and direct visual planning.

## II. RELATED WORK

We discuss prior approaches to tool use along with robotic learning methods that use demonstrations or video prediction.

**Planning tool use with analytic models.** Robotic manipulation involving tools has been studied in the task and motion planning (TAMP) literature [24, 25, 19, 48, 18, 32]. [45, 8] propose to use logic programming together with known models to algorithmically discover tool-use. One challenge that limits the scalability of most logic-based systems and analytic model-based systems is that modeling errors quickly accumulate during execution, which often results in fragile system. Unlike this prior work, we study tool use in the real world with visual inputs, using learned dynamics models and sampling-based optimization.

**Direct learning methods for tool use.** Several works have decomposed tool use into multiple stages [8, 44]: tool selection, task-oriented grasping of the selected tool [31, 41, 5, 2, 14], and using the grasped tool through planning [29] or policy learning [14]. These methods constrain the scope of motions to trajectories that involve the tool, while our method is capable of finding plans with or without a tool based on the situation (see, e.g., Figure 13). Other approaches have learned dynamics models to predict the outcome of applying actions to a tool [42, 29]. Unlike these approaches, which either use hand-designed perception pipelines or no visual perception at all, our method learns about object interactions directly from

raw image pixels, avoiding restrictive assumptions that might hinder generalization.

**Learning from demonstrations.** Imitating expert demonstrations is a common approach for learning complex skills and can overcome the exploration challenges arising in long-horizon control problems [3]. Prior work has leveraged demonstrations to accelerate *model-free reinforcement learning* either in simulation or the real world, overcoming the well-known exploration problem [38, 33, 21, 22]. Unlike most of these works, our method does not fully rely on a policy to obtain actions. Further, our method can leverage demonstrations of many different tasks and goals, by combining a stochastic policy learned from imitation with goal-directed model-based planning. Learning from demonstrations has also been used in combination with planning , where a planning cost function is inferred from data [39, 49, 1, 26] or where tool-use is learned from human demonstrations using pose-tracking [50, 28]. These approaches can be brittle since accurate pose-tracking of objects is often challenging in real-world scenarios.

**Video prediction-based planning** Our work extends visual MPC [16, 13], also referred to as visual foresight, which is a model-based reinforcement learning approach where a deep neural network model is trained to predict future visual observations. Such methods have been used in prior work for reaching [9], pushing objects [16, 11], basic grasping and relocation [12], and manipulating clothes [13]. Our aim is to extend these methods to enable improvisational tool use. Visual MPC, as well as other video prediction-based planners [9, 16, 36, 46], generally do not succeed at such temporally extended tasks. To this end, we propose to incorporate demonstrations into the algorithm to enable multi-stage tool-use capabilities, while still retaining the flexibility of goal-directed planning to accomplish varied user-specified goals.

## III. CAPABILITIES FOR IMPROVISATIONAL TOOL USE

Our goal is to study how robots can use novel, previously-unseen objects as tools in order to perform tasks that cannot be completed without tools. See Figure 3 for an example: the goal is for the robot to move the clutter onto the dustpan. Despite the robot having never previously seen the scraper, the clutter, or the
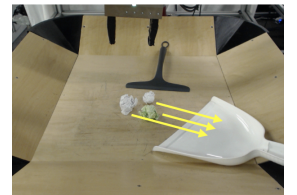


Fig. 3. How can a robot use a tool that it has never seen before to accomplish the goal?

dustpan in the scene, we would like the robot to figure out how it might use the scraper as a tool to efficiently clean the clutter. We hypothesize that one way to accomplish this is for the robot to learn to predict the consequences of its actions and the outcomes of object-object interactions. That is, to learn *generally* about how different objects interact with each other. Building realistic models of multi-object, non-prehensile interactions is challenging with analytic methods, as friction and contact dynamics become extraordinarily complex [15], and inferring them directly from images is
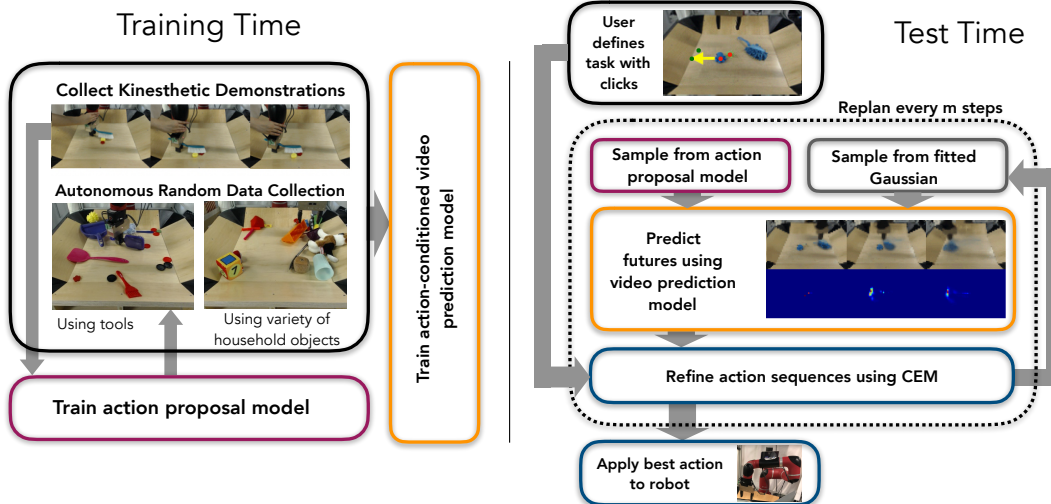
Fig. 2. Our guided visual foresight (GVF) approach, at training time (left) and test time (right). Our method incorporates demonstrations and autonomous data collection to learn a video prediction model and action proposal model that enable the robot to solve both a diverse range of goals that require tool use. We incorporate the action proposal model both for training data for the video prediction model and for improving the sampling-based planner at test time. The test time procedure is further detailed in Algorithm 1.

extremely difficult. Further, such models are not extendable to previously unseen objects, without requiring detailed physical knowledge of the object. In contrast, it is relatively easy for a robot to autonomously collect substantial amounts of data of object-object interactions. Thus, if we are able to *learn* a model from such interaction data, we can then use such a model for planning to use tools.

However, even with a learned model, generalization remains a critical challenge. In order for a robot to be able to plan to use novel, previously-unseen objects as tools, the robot needs a representation that can effectively generalize to new characteristics of objects such as new sizes, shapes, and masses. How then should we represent the objects and the environment? One option is to represent objects and surfaces using 3D meshes or voxel grids. However, this puts significant stress on acquiring a robust perception system that generalizes to novel tools, and would likely require significant supervised or simulation data. An alternative option is to use raw sensor readings, such as image pixels, as the representation. While such a representation does not incorporate object-centric inductive biases, it does have a number of benefits. First, we can train models of sensory observations, i.e. video prediction models, from completely unlabeled interaction data requiring no manual annotation. Second, low-level sensory observations such as pixels include all information about the environment that the robot can currently perceive, and hence are general to a wide range of objects and situations, including nonrigid and deformable objects. Motivated by these benefits, we will explore how we might enable improvisational tool use by autonomously collecting data of diverse object interactions, training predictive models of low-level sensory observations (i.e. action-conditioned video prediction [17]), and using these models to make plans to achieve goals involving tools. In the next section, we will describe how we can extend visual MPC [16, 13] to allow us to study such complex tasks.

## IV. DEMONSTRATION-GUIDED VISUAL PLANNING

We aim to use demonstrations to better enable visual MPC, or visual foresight, to perform more complex, temporally-extended tasks. While demonstrations are typically used with single-task imitation learning, we hope to incorporate the demonstrations in a way that maintains the generality of visual MPC. That is, we want both breadth and depth: a method that can be used both for solving a *variety* of tasks with unseen objects, and for solving a variety of *complex* goals, such as picking up a tool and using it. To this end, we will collect demonstrations that cover a broad range of tasks and goals using a range of tools. These demonstrations will be used in two ways: for improving the video prediction model and for improving the sampling-based planning process. Specifically, we will use demonstrations, first, to enable the robot to collect data in parts of the state space that the robot would be unlikely to visit with random interaction, and second, to aid the sampling-based planner in finding solutions that are more difficult to find when searching from scratch. We next give a short overview of our guided visual foresight approach (GVF), and then describe each of the components in more detail.

As shown in the left of Figure 2 and described in subsection IV-A, the *training time* procedure consists of three parts. We collect kinesthetic demonstrations of trajectories involving tool use. We use this data to train an *action proposal model* to obtain a distribution over action sequences conditioned on the initial image based on actions taken by the demonstrator. This proposal model will be used both during training, for collecting data for improving the video prediction model, and at test time, to help warm-start the optimization over action sequences. In addition to the human demonstrations, the robot autonomously collects data by executing random actions. Finally, we train an action-conditioned video prediction model to predict future video sequences based on the initial image and the corresponding action sequences.
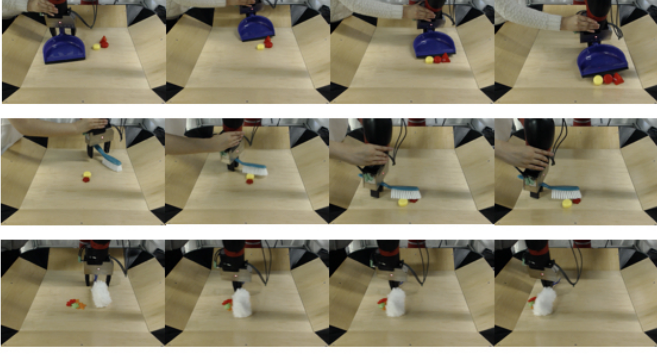
Fig. 4. Examples of demonstrations collected via kinesthetic teaching. We provide demonstrations for a diverse range of tools, objects, and goals.

As illustrated on the right in Figure 2 and described in subsection IV-B, at *test time*, we enable the robot to plan to use objects as tools as follows. First, a user can specify a goal by clicking on objects in the image and selecting where the corresponding pixels should move. For example, the user might specify that three pieces of trash need to be moved to a location within a dustpan. Then, the current observation is passed to the action proposal model, which returns a sampling distribution that is used to sample a certain number of action sequences. These samples will usually correspond to different ways of interacting with objects in the scene — in our case different ways of using objects in the scene as tools. We feed each of the sampled action sequences into the video prediction model to predict their outcome as a video. We then rank these predictions using a cost function determined by the human-specified goal, and refine the best samples further. Lastly, the robot recomputes action plans after several control cycles.

Because the demonstrations entail a wide variety of tools and goals, our experiments find that pure imitation learning struggles to capture the breadth of such a distribution. However, these demonstrations can be effectively used to guide the planning process towards tool-related behaviors, while the predictive model is used to fully construct and refine a sequence of actions for completing the task.

*A. Training Time*

*1) Demonstration data collection:* We collect demonstrations of tasks that require tools, typically involving grasping a tool and using the tool to maneuver other objects to a certain location. Because we specifically care not just about accomplishing a single task, but being able to perform a range of tasks with many different objects, we collect demonstrations for a variety of tasks that require a variety of tools (instead of many demonstrations for a single task). In our prototype, we use kinesthetic teaching to collect demonstrations. During the demonstrations we record images $I_t$, Cartesian end-effector positions $s_t$, as well as motor commands $a_t$ in the form of Cartesian end-effector displacements. We will denote the demonstration data as $\mathcal{D}_{\text{demo}} = \{(I_1, s_1, a_1, I_2, s_2, a_2, ...)_j\}$ For examples of demonstrations, see Figure 4. For each demonstration, we record a sequence of 24 to 30 time-steps.

The tasks for these demonstrations are chosen such that success without tool use would be very low.

*2) Action proposal model training:* In order to collect additional data of interaction with tools and guide the planning process at test-time, we aim to acquire an approximate model of the tool-use behaviors seen in the demonstrations. We fit an action proposal model $g_\theta$ to the demonstration data that outputs the distribution over a sequence of actions conditioned on the initial image $I_1$ and robot joint positions $s_1$. Note that outputting a distribution, rather than a deterministic output, will allow the model to capture a range of behaviors present in the diverse demonstrations. Because we would like to use this model for sampling-based planning, we do not condition the model on the final image of the demonstration. The action proposal model is parameterized as an autoregressive recurrent neural network (RNN). It is trained with the following maximum likelihood objective:

$$\max_{\theta} \sum_{I_1, s_1, a_{1:T} \in \mathcal{D}_{\text{demo}}} \log p_\theta(a_1, ..., a_T | I_1, s_1) \qquad (1)$$

Here, $I_1, s_1, a_{1:T}$ are the initial image observation, initial end-effector position, and action sequence in a kinesthetic demonstration, $p_\theta$ is mixture of Gaussians where the parameters are produced by $g_\theta$. We use a long short-term memory network (LSTM) [20] to model recurrence. The model $g_\theta$ consists of three components: an initial state encoder $g_e$, an action encoder $g_a$, and an LSTM cell. The encoder $g_e$ encodes the initial image $I_1$ and state $s_1$ to provide the input to the LSTM at the first timestep, and the action encoder $g_a$ encodes the previous action to provide to the LSTM cell and future timesteps. The recurrent LSTM cell produces both the parameters of the Gaussian mixture and the next hidden state $h_{t+1}$. The full network is described by the following equations.

$$a_0 = \mathbf{0} \qquad (2)$$
$$h_0, \varnothing = \text{LSTM}(g_e(I_1, s_1), \mathbf{0}) \quad t = 0 \qquad (3)$$
$$h_t, (\mu_t^{(i)}, \Sigma_t^{(i)}, w_t^{(i)}) = \text{LSTM}(g_a(a_{t-1}), h_{t-1}) \ \forall t > 0 \quad (4)$$
$$p(a_t) = \sum_{i \in \{1,..,N_c\}} w_t^{(i)} \mathcal{N}(\hat{a}_t; \mu_t^{(i)}, \Sigma_t^{(i)}) \quad (5)$$

The zero action is used to indicate the start of the sequence.

We illustrate the neural network architecture for our action proposal model in Figure 5. To handle RGB image inputs, the network $g_e$ is composed of three convolutional layers with 32 3×3 filters, each followed by batch normalization and a ReLU non-linearity. The first convolutional layer is initialized with pretrained weights from VGG-16. Then, spatial feature points are extracted from the last convolutional layer with a spatial soft-argmax operation [30]. We concatenate these features with the robot's initial end-effector position $s_1$ and pass them through a fully-connected layer with 50 units. As previously discussed, we use an RNN, in particular an LSTM net, that uses this embedding to initialize the hidden state. In the action encoder $g_a$, the actions are passed through a fully-connected layer, the result of which are the inputs to the
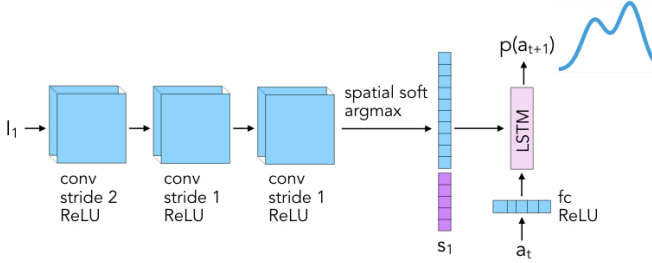
Fig. 5. Architecture for the action proposal model. We use a recurrent autoregressive model to output the parameters of a Gaussian mixture model over the probbaility of an action at each timestep. Using recurrence and Gaussian mixtures enables the network to model diverse and multi-modal demonstration data based on the initial image.



Fig. 6. Autonomous data collection with a collection of house-hold objects (left) and a variety of of tools (right)

LSTM. Finally, we use a mixture density network with 10 components to represent the output distribution [7].

*3) Autonomous data collection:* For training the predictive model, we need trajectory data $(I_1, s_1, a_1, I_2, s_2, a_2, ...)$ from the robot's interactions. The demonstration data, on its own, is quite limited in quantity. Hence, we also choose to have the robot collect data autonomously using, for example, random actions sampled from a Gaussian distribution. However, when training a video prediction model using this data, we observe that the model makes optimistic predictions, very often predicting that the robot will grasp a tool when approaching it, even when the grasp is not correctly positioned. To address this issue, the robot additionally collects data autonomously by sampling actions from the action proposal model.

All in all, the video prediction model is trained on three sources of data: demonstrations, random interactions, and on-policy data. When sampling the actions from the action proposal model a variety of tools are present in the robot workspace. When executing random actions, in some parts of the data we use tools, see Figure 6 (right), whereas in the other parts we use a collection of household items, shown in Figure 6 (left). Further details about the composition of the dataset are given in section V.

To simplify picking up objects in general, including tools, we incorporate a simple "grasp-reflex" into the controller, where the gripper automatically closes when the height of the wrist above the table is lower than a small threshold (following Ebert et al. [13]). This reflex is inspired by the palmar reflex observed in infants [40]. With this primitive, when collecting data with random actions and rigid objects, about 20% of trajectories included some sort of object grasp,

significantly higher than without the reflex. Note, however, that this technique on its own is not sufficient for enabling tool use, as we find in the experiments.

*4) Predictive model training:* Once we collect autonomous data, we use it to build a predictive model of future sensory inputs, i.e. images, conditioned on the initial image and the future actions taken. We use a transformation-based video prediction architecture, first proposed by Finn et al. [17], and use the open-source architecture from Ebert et al. [13]. The advantage of using transformation-based models over a model that directly generates pixels is two-fold: (1) prediction is easier, since the appearance of objects and the background scene can be reused from previous frames and (2) the transformations can be leveraged to obtain predictions about *where* pixels will move, a property that is used in our planning cost function formulation, presented in subsubsection IV-B1. The model, which is implemented as a recurrent convolutional neural network, $f_\gamma$ parameterized by $\gamma$, has a hidden state $h_t$ and takes in a previous image and an action at each step of the rollout. Future images $\hat{I}_{t+1}$ are generated by warping the previous generated image $\hat{I}_t$ or the previous true image $I_t$, when available, according to a 2-dimensional flow field $\hat{F}_{t+1\leftarrow t}$. The forward pass of the dynamics model is summarized in the following two equations:

$$[h_{t+1}, \hat{F}_{t+1\leftarrow t}] = f_\gamma(a_t, h_t, I_t) \qquad (6)$$

$$\hat{I}_{t+1} = \hat{F}_{t+1\leftarrow t} \diamond \hat{I}_t \qquad (7)$$

The model is trained with stochastic gradient descent using a $\ell_2$ image reconstruction loss. For more details on the architecture and training, see Appendix VII-A.

### B. Test-Time Control

At test time, a user provides a goal to the robot, and the robot uses the learned action proposal and video prediction models to plan to achieve the goal. We describe this process in more detail next.

*1) Planning cost function:* A user can provide a goal by clicking on a pixel corresponding to an object and a corresponding goal position for that pixel. A *pixel distance cost function* evaluates how far the designated pixel is from the goal pixels. Given a distribution over pixel positions $P_0$, our model predicts distributions over its positions $P_t$ at time $t \in \{1, ..., T\}$ as follows: To predict the future positions of the designated pixel $d$, the same transformations used to transform the images are applied to the distribution over designated pixel locations. The warping transformation $\hat{F}_{t+1\leftarrow t}$ can be interpreted as a stochastic transition operator allowing us to make probabilistic predictions about future locations of individual pixels:

$$\hat{P}_{t+1} = \hat{F}_{t+1\leftarrow t} \diamond \hat{P}_t \qquad (8)$$

Here, $P_t$ is a distribution over image locations which has the same spatial dimension as the image (an example is shown in Figure 8 in the third row). For simplicity in notation, we will use a single designated pixel moving forward, but using

**Algorithm 1** Guided visual foresight (test time)

1: **Input:** Predictive model $f_\gamma$
2: **Input**: Planning cost $c$ derived from user-specified pixel goals
3: **for** $i = 0...n_{iter} - 1$ **do**
4:     **if** $i == 0$ **then**
5:         Sample $M$ action sequences $\{a_{1:H}^{(m)}\}$ from action proposal distribution by rolling out $g_\theta$
6:     **else**
7:         Sample $M$ action sequences $a_{1:H}^{(m)}$ from $\mathcal{N}(\mu^{(i)}, \Sigma^{(i)})$
8:     Use $f_\gamma$ to predict future image sequences $\hat{I}_{1:H}^{(m)}$ and probability distributions $\hat{P}_{1:H}^{(m)}$
9:     Rank action sequences using cost function $c$
10:    Fit a Gaussian to the $k$ action samples with lowest cost, yielding $\mu^{(i+1)}, \Sigma^{(i+1)}$
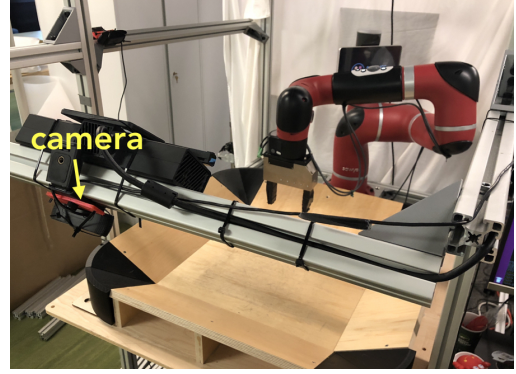11: Execute lowest-cost action sequence on the robot

---



Fig. 7. The physical robot set-up: we use a Sawyer robot with Cartesian space impedance control to ensure soft interaction with objects. RGB images are taken from a conventional webcam, as indicated in the image.

multiple designated pixels is a straightforward extension. At the first time step, the distribution $\hat{P}_0$ is defined to be 1 at the position of the user-selected designated pixel and zero elsewhere. One way of defining the cost per time-step $c_t$ is by using the expected Euclidean distance to the goal point $d_g$, which is straight-forward to calculate from $P_t$ and $g$, as follows:

$$c = \sum_{t=1,...,T} c_t = \sum_{t=1,...,T} \mathbb{E}_{\hat{d}_t \sim P_t} \left[ \| \hat{d}_t - d_g \|_2 \right] \qquad (9)$$

The per time-step costs $c_t$ are summed together giving the overall planing objective $c$. For tasks with multiple designated pixels $d^{(i)}$, the costs are also summed together.

*2) Planning with demonstration guidance:* Planning with GVF at test time is illustrated in Figure 2 (right) and Algorithm 1. The user first specifies the task by clicking on the pixels that shall be moved and the corresponding goal-pixels. The planner searches for actions using the cross entropy method (CEM) [23], a common iterative sampling-based optimization procedure. To allow the optimizer to find more complicated, temporally extended action sequences, such as picking up a tool and using it in a goal-directed manner, we sample actions from the stochastic action proposal model $g_\theta$ in the first iteration of CEM, as listed in line 5 of Algorithm 1. After rolling out the video prediction model $f_\gamma$ using Equation 7 we obtain $M$ different predicted probability distributions $\hat{P}_{1:H}^m$, which are ranked using the cost function $c$. We then fit a Gaussian distribution to the best $k$ action samples (see line 10). In later CEM iterations, actions are sampled from the fitted Gaussians (line 7). In practice, we choose the number of samples $M$ to be 100, the horizon $H$ to be 21, and the number of CEM iterations $n_{iter}$ to be 3.

## V. EXPERIMENTS

In our experiments, we aim to answer the following questions: **(1)** can our approach effectively solve tasks that require tool use? **(2)** can our method improvise, by figuring out how

to use a new object that was not seen during training as a tool? **(3)** how important is the action proposal model? **(4)** can our method dynamically decide to use or not use tools depending on the demands of the task? To answer these questions, we conduct experiments on a Sawyer robotic arm, with an experimental set-up shown in Figure 7. Video results are in the supplemental materials and the supplementary webpage[2].

### A. Experimental Set-Up and Comparisons

To train the action proposal model discussed in Section IV-A2, we collected kinesthetic demonstrations on a Sawyer robotic arm with twenty different tools. Figure 9 (left) illustrates these twenty tools. Here, we focus on sweeping, scraping, and wiping tasks where the goal is to move multiple objects which would be infeasible to complete without the use of a tool. Thus, in each demonstration, we randomly place a tool and pile of objects in front of the robot and kinesthetically demonstrate how to grasp the tool and sweep the smaller objects (see Figure 4). The demonstrations are recorded at 5 Hz and range from 24 to 30 time steps. The action proposal model is then trained on subsequences of 10 steps, conditioned on the image observation and robot state of the first step of the subsequence.

As described in Section IV-A4, to train the video prediction model, we collect additional interaction data by taking random actions and by rolling out samples from our action proposal model. Our final dataset is composed of: 16,000 random trajectories from the open-source dataset in [12], 5,052 random trajectories with tools, 1,754 samples from the action proposal model, and 1,200 demonstrations.

To study the importance of visual planning and the importance of demonstrations in both video prediction model training and planning, we compare our method to the following approaches:

- **Imitation Learning**: Sampling from an action proposal model that is conditioned on the initial and goal image observations, representative of standard imitation learning [37, 4, 10, 35]. This comparison evaluates the importance of physical prediction and planning.

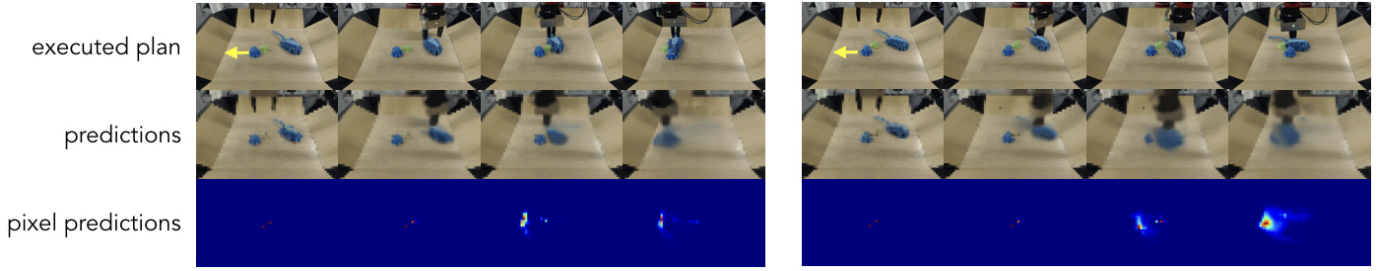[2]For videos, see https://sites.google.com/view/gvf-tool

Fig. 8. Examples of the lowest-cost predictions (2nd row) and executed actions (1st row), for the task indicated in the top left image. The left example shows a model trained on on-policy data, while the right example shows the best action sequence found with a model that was *not* trained on demonstration data nor on data from the action proposal model. Note that the robot fails to grasp the object in the second example, while the model predicted that the grasp would be successful. Each example also shows the probability distribution of the designated pixel over time (3rd row).



Training Tools          Test Tools

Fig. 9. Left: tools used during training. Right: test tools used in our quantitative evaluation, some of which are not conventional tools.
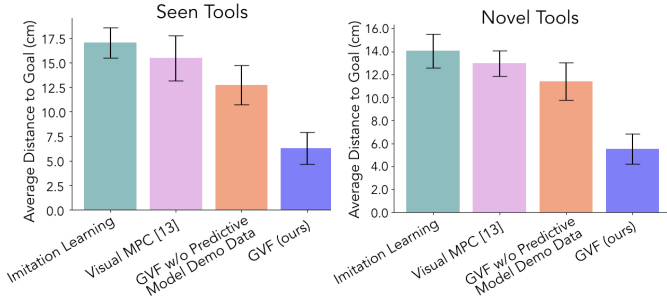


Fig. 10. Quantitative results: our approach, which uses elements of imitation learning and visual MPC, significantly outperforms either approach individually. In particular, we compare to direct imitation learning on the diverse demonstration data, to visual MPC without the learned action proposal model, and to our method with a video prediction model trained only on autonomously collected data.

- **Visual MPC**: Our method with CEM samples from a Gaussian distribution. This comparison is representative of visual MPC [16, 13]. The video prediction model is still trained with demonstrations and samples from the action proposal model, so it is actually stronger than the method of Ebert et al. [13]. This comparison evaluates the importance of demonstrations in guiding the planning process.
- **GVF w/o Predictive Model Demo Data**: Our method with a video prediction model trained only on randomly collected data, omitting demonstration data and data from the action proposal model when training the video prediction model. Test-time planning still uses the action proposal model. This comparison evaluates the role of demonstrations in improving the predictions of multi-object interactions.

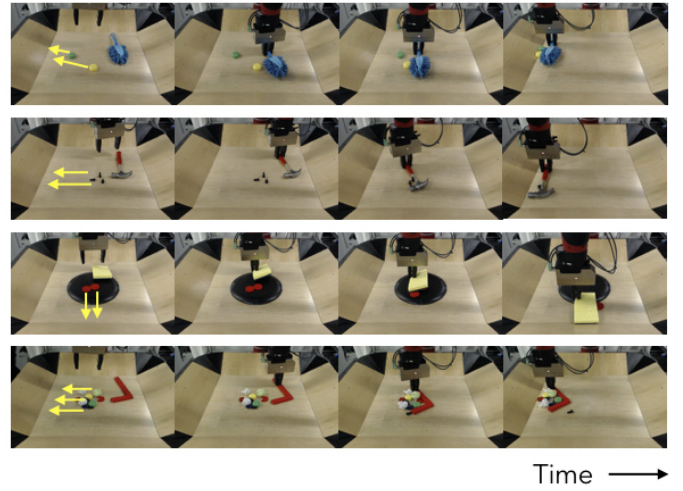The imitation learning method uses a convolutional neural



Fig. 11. Qualitative results illustrating our approach (GVF): the robot executing the lowest-cost plan according to the pixel-distance cost function. The arrows in the left column of images indicate where the robot is to maneuver the objects.

network with the same architecture as for the action proposal model, described in Section IV-A2 to map both the initial and goal image observations. The feature points from both image inputs, along with the robot's initial end-effector position, are concatenated and passed through a fully-connected layer of 50 units. The output is then used as the initial state of the LSTM.

### B. Experimental Results

We quantitatively evaluate each method on 10 tasks with tools seen during training and 10 tasks with previously unseen tools, with results summarized in Figure 10 and detailed in Appendix VII-B. Each task requires picking up the tool and sweeping, scraping, or wiping objects to the position corresponding to the specified goal pixels. Note that the set of tasks with seen tools differs from the set of tasks with novel tasks, thus the two sets of results are not directly comparable. In regard to question **(1)**, these results show that our method can successfully use tools when they are available, reaching less than half the position error to the goal compared to the prior methods. The qualitative results and supplementary video show that the robot is generally successful at the tool

GVF (ours)

GVF w/o Predictive Model Demo Data

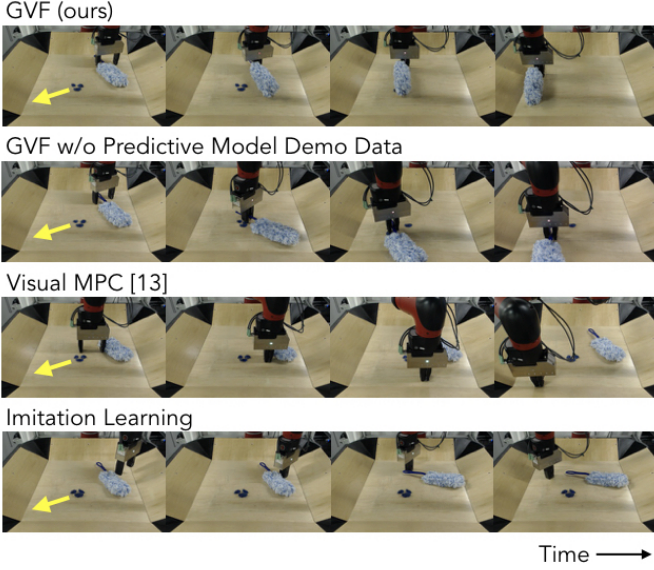Visual MPC [13]

Imitation Learning

Time ⟶

Fig. 12. Qualitative comparison of our method to prior methods and ablations. Without the action proposal model (visual MPC), the method generally cannot find the actions that grasp the tool, while without the video prediction model, imitation learning generally fails to use the tool in a meaningful way.
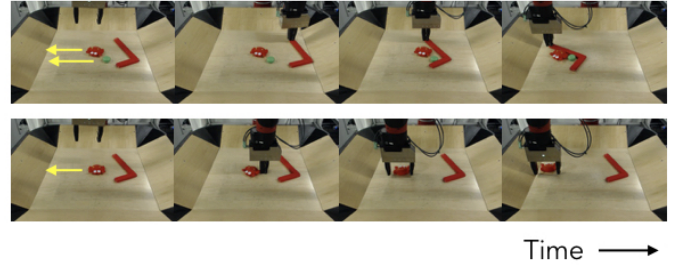


Time ⟶

Fig. 13. Our approach solves a task using a tool where it needs to manipulate two objects simultaneously (top), and chooses to not use a tool when the task involves only one object, allowing the robot to complete the task more efficiently (bottom).

use tasks. In regard to question **(2)**, the relatively similar performance with novel tools indicates that the robot can generalize effectively, utilizing new objects as tools when needed. The qualitative results illustrate examples of novel objects that the robot can use as tools.

We further evaluate questions **(1)** and **(2)** through qualitative results and comparisons. Our primary qualitative results are in Figures 1 and 11, including both seen and novel tools. As illustrated in Figure 11 (third row), we find that, even though the robot has never seen a sponge before, the robot is able to use the sponge to wipe a plate. Further, as shown in Figure 1 (second row), the robot has never seen or interacted with a knife before, it can figure out how to use the knife to push two red pieces of rubbish to the edge of the cutting board. Finally, when no conventional tool is available, the robot is able to improvise when tasked with moving pieces of trash, by grasping a water bottle and using it to sweep the trash to the side (see Figure 1, third row).

In regard to question **(3)**, incorporating demonstration data for both the sampling distribution and training the video prediction model on average leads to more successful behavior. Moreover, compared to learning a predictive model with demonstration data, the action proposal model results in better performance on our set of evaluation tasks. We analyze the failure modes of the prior methods and ablations in Figure 12.

Lastly, in regard to question **(4)**, we aim to test whether our method still retains the full *generality* of visual MPC by evaluating whether it can plan to solve a task *without* using a tool, e.g. when tool-use does not have an advantage, or even a disadvantage. Therefore we set up the following experiment: as shown in Figure 13, we have two almost identical task settings, where the only difference is that in the first task two objects need to be pushed and in the second task only one object needs

to be moved. In the latter case using a tool does not have any advantage, therefore we expect the planner to find a plan that does not use the tool. To allow the planner to explore both tool-use and non-tool use options, in the first CEM-iteration we use a combination of 50 samples from the proposal model and 50 samples from a unit Gaussian. As shown in Figure 13, GVF is indeed able to find a non-tool use trajectory for the single-object pushing task.

## VI. DISCUSSION

**Summary.** We developed an approach to enable a robot to accomplish both *diverse* and *complex* tasks involving previously-unseen objects with access to only raw visual inputs. We studied the particular case of solving many different tasks that require manipulating objects as tools. Our approach learns from a combination of diverse human demonstration data, with many different goals, tools, and items, and autonomously-collected interaction data, with diverse objects. We show how we can use this data to train a model that can predict the visual outcome of actions that cause multi-object interaction, and use these predictions to figure out how to accomplish tasks by leveraging such object-object interactions.

**Limitations and Future Work.** Our approach has a number of limitations that we hope to study in future work. First, the tool-use tasks that we consider are diverse, but largely involve sweeping, wiping, and hooking interactions. In future work, we hope to also study tool use problems that involve cutting, skewering, and screwing interactions. In these cases, we expect that a more unconstrained action space may be important, where demonstrations may be of even greater importance to direct exploration within the larger state space. Second, our approach uses entirely visual observations, while a number of tool-use applications, such as using a screwdriver, demand force feedback. In principle, our approach abstracts away the form of the observation through learning. Indeed, prior work has shown that approaches like visual foresight can be integrated with tactile sensor inputs [43]; but in practice, the introduction of tactile or force sensors would likely introduce additional challenges in evaluating predictions and collecting data safely. Finally, while tool use provides an interesting testbed for studying diverse, yet complex manipulation problems, we hope to study our approach in the context of other temporally-extended skills in future work.

REFERENCES

[1] Jacopo Aleotti and Stefano Caselli. Grasp recognition in virtual reality for robot pregrasp planning by demonstration. In *International Conference on Robotics and Automation (ICRA)*, 2006.

[2] Rika Antonova, Mia Kokic, Johannes A Stork, and Danica Kragic. Global search with bernoulli alternation kernel for task-oriented grasping informed by simulation. In *Conference on Robot Learning (CoRL)*, 2018.

[3] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Journal of Robotics and Autonomous Systems*, 2009.

[4] Christopher G Atkeson and Stefan Schaal. Robot learning from demonstration. In *International Conference on Machine Learning (ICML)*. Citeseer, 1997.

[5] Yasemin Bekiroglu, Dan Song, Lu Wang, and Danica Kragic. A probabilistic framework for task-oriented grasp stability assessment. In *International Conference on Robotics and Automation (ICRA)*, 2013.

[6] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[7] Christopher M Bishop. Mixture density networks. Technical report, Citeseer, 1994.

[8] Solly Brown and Claude Sammut. A relational approach to tool-use learning in robots. In *International Conference on Inductive Logic Programming*. Springer, 2012.

[9] Arunkumar Byravan and Dieter Fox. Se3-nets: Learning rigid body motion using deep neural networks. In *International Conference on Robotics and Automation (ICRA)*, 2017.

[10] Felipe Codevilla, Matthias Miiller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2018.

[11] Frederik Ebert, Chelsea Finn, Alex X Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. *Conference on Robot Learning (CoRL)*, 2017.

[12] Frederik Ebert, Sudeep Dasari, Alex X Lee, Sergey Levine, and Chelsea Finn. Robustness via retrying: Closed-loop robotic manipulation with self-supervised learning. *Conference on Robot Learning (CoRL)*, 2018.

[13] Frederik Ebert, Chelsea Finn, Sudeep Dasari, Annie Xie, Alex Lee, and Sergey Levine. Visual foresight: Model-based deep reinforcement learning for vision-based robotic control. *arXiv preprint arXiv:1812.00568*, 2018.

[14] Kuan Fang, Yuke Zhu, Animesh Garg, Andrey Kurenkov, Viraj Mehta, Li Fei-Fei, and Silvio Savarese. Learning task-oriented grasping for tool manipulation from simulated self-supervision. *Robotics: Science and Systems (RSS)*, 2018.

[15] Nima Fazeli, Elliott Donlon, Evan Drumwright, and Alberto Rodriguez. Empirical evaluation of common contact models for planar impact. In *International Conference on Robotics and Automation (ICRA)*, 2017.

[16] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *International Conference on Robotics and Automation (ICRA)*, 2017.

[17] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.

[18] SK Gupta, Christiaan JJ Paredis, and PF Brown. Micro planning for mechanical assembly operations. In *International Conference on Robotics and Automation (ICRA)*, volume 1. IEEE, 1998.

[19] Dan Halperin, J-C Latombe, and Randall H Wilson. A general framework for assembly planning: The motion space approach. *Algorithmica*, 2000.

[20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997.

[21] Jens Kober, Erhan Öztop, and Jan Peters. Reinforcement learning to adjust robot movements to new situations. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2011.

[22] Petar Kormushev, Sylvain Calinon, and Darwin G Caldwell. Robot motor skill coordination with em-based reinforcement learning. In *International Conference on Intelligent Robots and Systems (IROS)*, 2010.

[23] Dirk P Kroese, Reuven Y Rubinstein, Izack Cohen, Sergey Porotsky, and Thomas Taimre. Cross-entropy method. In *Encyclopedia of Operations Research and Management Science*. 2013.

[24] James Kuffner, Koichi Nishiwaki, Satoshi Kagami, Masayuki Inaba, and Hirochika Inoue. Motion planning for humanoid robots. In *Robotics Research. The Eleventh International Symposium*. Springer, 2005.

[25] Jean-Claude Latombe. *Robot motion planning*, volume 124. Springer Science & Business Media, 2012.

[26] Martin Lawitzky, Jose Ramon Medina, Dongheui Lee, and Sandra Hirche. Feedback motion planning and learning from demonstration in physical robotic assistance: differences and synergies. In *International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2012.

[27] Alex X Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine. Stochastic adversarial video prediction. *arXiv:1804.01523*, 2018.

[28] Dongheui Lee, Hirotoshi Kunori, and Yoshihiko Nakamura. Association of whole body motion from tool

knowledge for humanoid robots. In *International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2008.

[29] Ian Lenz, Ross A Knepper, and Ashutosh Saxena. Deepmpc: Learning deep latent features for model predictive control. In *Robotics: Science and Systems (RSS)*, 2015.

[30] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research (JMLR)*, 2016.

[31] Zexiang Li and S Shankar Sastry. Task-oriented optimal grasping by multifingered robot hands. *Journal on Robotics and Automation*, 1988.

[32] Igor Mordatch, Emanuel Todorov, and Zoran Popović. Discovery of complex behaviors through contact-invariant optimization. *ACM Transactions on Graphics (TOG)*, 2012.

[33] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Overcoming exploration in reinforcement learning with demonstrations. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2018.

[34] François Osiurak and Arnaud Badets. Tool use and affordance: Manipulation-based versus reasoning-based approaches. *Psychological review*, 2016.

[35] Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Dian Chen, Yide Shentu, Evan Shelhamer, Jitendra Malik, Alexei A Efros, and Trevor Darrell. Zero-shot visual imitation. In *International Conference on Learning Representations (ICLR)*, 2018.

[36] Chris Paxton, Yotam Barnoy, Kapil Katyal, Raman Arora, and Gregory D Hager. Visual robot task planning. *arXiv preprint arXiv:1804.00062*, 2018.

[37] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1989.

[38] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *Robotics: Science and Systems (RSS)*, 2018.

[39] N. Rhinehart, R. McAllister, and S. Levine. Deep Imitative Models for Flexible Inference, Planning, and Control. *ArXiv e-prints*, October 2018.

[40] D. Sherer. Fetal grasping at 16 weeks' gestation. *Journal of ultrasound in medicine*, 1993.

[41] Karun B Shimoga. Robot grasp synthesis algorithms: A survey. *International Journal of Robotics Research (IJRR)*, 1996.

[42] Alexander Stoytchev. Behavior-grounded representation of tool affordances. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2005.

[43] Stephen Tian, Frederik Ebert, Dinesh Jayaraman, Mayur Mudigonda, Chelsea Finn, Roberto Calandra, and Sergey Levine. Manipulation by feel: Touch-based control with deep predictive models. *International Conference on Robotics and Automation (ICRA)*, 2019.

[44] V Tikhanoff, U Pattacini, L Natale, and G Metta. Exploring affordances and tool use on the icub. In *International Conference on Humanoid Robots (Humanoids)*, 2013.

[45] Marc Toussaint, Kelsey Allen, Kevin Smith, and Joshua B Tenenbaum. Differentiable physics and stable modes for tool-use and manipulation planning. *Robotics: Science and Systems (RSS)*, 2018.

[46] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[47] SHI Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[48] Atsushi Yamashita, Jun Sasaki, Jun Ota, and Tamio Arai. Cooperative manipulation of objects by multiple mobile robots with tools. In *Proceedings of the 4th Japan-France/2nd Asia-Europe Congress on Mechatronics*, 1998.

[49] Gu Ye and Ron Alterovitz. guided motion planning. In *Robotics Research*. Springer, 2017.

[50] Yixin Zhu, Yibiao Zhao, and Song Chun Zhu. Understanding tools: Task-oriented object modeling, learning and recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

## A. Video Prediction Model Implementation details

We deterministic variant of the video prediction model described in [27]. The video prediction model is realized as a transformation-based model which generates future images by *transforming* past images. The core of the model is made up from a recurrent neural network, Figure 14 gives an overview of a roll-out through time. In practice, the first two images passed into the model are ground truth images, called *context frames*. At every time-step an action $a_t$ is passed into the model along with the hidden state $h_t$, producing a new state $h_{t+1}$ and a flow field $\hat{F}_{t+1 \leftarrow t}$ which is used to to transform the image via bi-linear sampling.

Figure 15 shows the forward pass of a single time-step. The network consists of multiple layers of convolutional LSTMs [47], a spatial, convolutional version of standard LSTMS, which are more efficient computationally and provide a regularizing inductive bias. While the transformations in theory would be sufficient to predict most parts of a video, it was found that allowing the model to selectively copy parts of the image from the *first frame* of the sequence helps overcoming problems that occur with *occluding objects*, i.e. objects in the fore-ground, would erase other parts of the image when they are moving [11]. Copying parts from the first image is achieved by predicting compositing masks (shown in green), a set of features maps with the same size of the image passed through a channel-wise softmax so that they add up to one along the channel-dimension.

The prediction loss is implemented as a standard $l1$-error. To regularize the RNN scheduled sampling [6] is used, which provides a training curriculum for more stable RNN training. The model is trained for 300k steps with standard back-propagation through time (BPTT), for the optimizer we use Adam. For more details concerning the video prediction model implementation we refer the reader to [13] and [27].
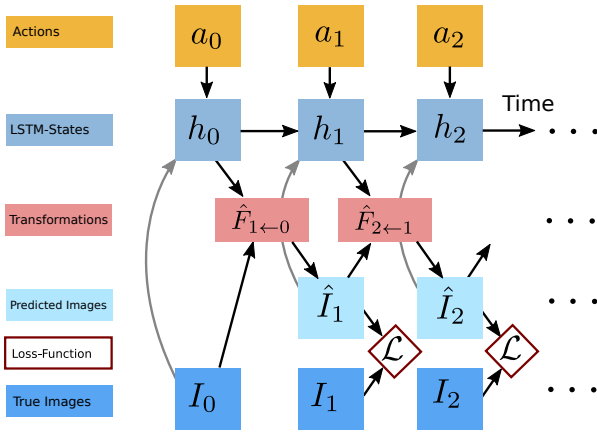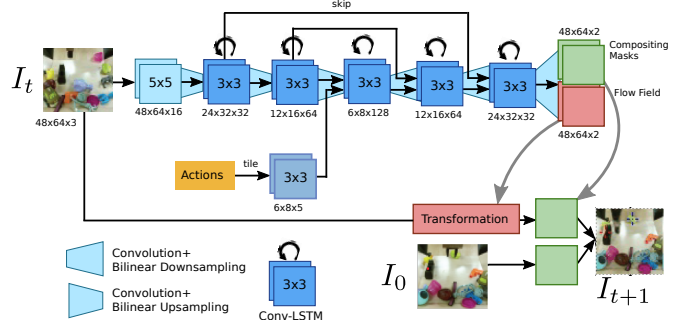


Fig. 15. Forward pass through the recurrent SNA model. The image from the first time step $I_0$ is concatenated with the transformed images $\hat{F}_{t+1 \leftarrow t} \diamond \hat{I}_t$ multiplying each channel with a separate mask to produce the predicted frame for step $t+1$. Used with permission from Ebert et al. [13].

## B. Detailed Quantitative Task Results

In Tables VII-B and VII-B, we show the complete results for all 20 evaluated tasks.

| Task | GVF (ours) | GVF w/o Predictive Model Demo Data | Visual MPC [13] | Imitation Learning |
|---|---|---|---|---|
| 1 | 10.0 | 20.0 | 20.0 | 18.5 |
| 2 | 8.5 | 7.0 | 8.5 | 14.5 |
| 3 | 0.5 | 13.0 | 15.7 | 15.7 |
| 4 | 9.5 | 5.5 | 19.3 | 18.5 |
| 5 | 1.0 | 16.7 | 24.3 | 24.3 |
| 6 | 0.3 | 8.0 | 7.5 | 14.5 |
| 7 | 5.0 | 18.3 | 22.0 | 18.3 |
| 8 | 13.5 | 18.7 | 18.7 | 20.0 |
| 9 | 13.7 | 18.7 | 18.7 | 21.0 |
| 10 | 1.7 | 1.7 | 0.5 | 5.0 |

Fig. 16. Average distance to goal (in centimeters) for each evaluation task with previously seen tools.

| Task | GVF (ours) | GVF w/o Predictive Model Demo Data | Visual MPC [13] | Imitation Learning |
|---|---|---|---|---|
| 1 | 3.0 | 12.8 | 11.5 | 14.0 |
| 2 | 8.3 | 10.0 | 15.5 | 10.0 |
| 3 | 2.0 | 0.0 | 9.5 | 9.0 |
| 4 | 1.5 | 19.5 | 18.0 | 21.5 |
| 5 | 4.5 | 11.0 | 12.0 | 11.5 |
| 6 | 12.5 | 19.3 | 19.3 | 18.8 |
| 7 | 3.0 | 9.5 | 8.0 | 13.0 |
| 8 | 0.3 | 11.0 | 11.5 | 11.0 |
| 9 | 10.5 | 11.0 | 14.5 | 22.0 |
| 10 | 10.0 | 10.0 | 10.0 | 10.0 |

Fig. 17. Average distance to goal (in centimeters) for each evaluation task with novel tools.



Fig. 14. Computation graph of the video prediction model. Time goes from left to right, $a_t$ are the actions, $h_t$ are the recurrent hidden states, $\hat{F}_{t+1 \leftarrow t}$ is a 2D-warping field, $I_t$ are real images, and $\hat{I}_t$ are predicted images, $\mathcal{L}$ is a pairwise training-loss. In this illustration $I_0$ is a context frame. Used with permission from Ebert et al. [13].