# arXiv:2109.14078v1 [cs.RO] 28 Sep 2021

# Learning Periodic Tasks from Human Demonstrations

Jingyun Yang<sup>1†</sup>, Junwu Zhang<sup>2</sup>, Connor Settle<sup>3</sup>, Akshara Rai<sup>4</sup>, Rika Antonova<sup>3‡</sup>, Jeannette Bohg<sup>3</sup>

Abstract-We develop a method for learning periodic tasks from visual demonstrations. The core idea is to leverage periodicity in the policy structure to model periodic aspects of the tasks. We use active learning to optimize parameters of rhythmic dynamic movement primitives (rDMPs) and propose an objective to maximize the similarity between the motion of objects manipulated by the robot and the desired motion in human video demonstrations. We consider tasks with deformable objects and granular matter whose state is challenging to represent and track: wiping surfaces with a cloth, winding cables/wires, stirring granular matter with a spoon. Our method does not require tracking markers or manual annotations. The initial training data consists of 10-minute videos of random unpaired interactions with objects by the robot and human. We use these for unsupervised learning of a keypoint model to get task-agnostic visual correspondences. Then, we use Bayesian optimization to optimize rDMPs from a single human video demonstration within few robot trials. We present simulation and hardware experiments to validate our approach.

# I. INTRODUCTION

Periodic tasks such as wiping a table with a cloth, stirring food, winding cables, or tying ropes are ubiquitous in our daily lives (see Figure 1). In this work, we address how robots can appropriately represent and learn periodic policies by watching humans. While prior works considered learning manipulation skills from human demonstrations [1]–[6], less attention has been given to periodic tasks. These tasks repeat similar motion with only small differences between repetitions. If a robot was able to decompose the demonstration of a periodic task into periods, it could efficiently learn the underlying motion and repeat it as many times as necessary. Prior hierarchical learning approaches investigated learning of compositional tasks either from demonstrations [7], [8] or through reinforcement learning [9], [10]. However, these approaches do not leverage the strong relationship between repetitions and can be inefficient for learning periodic tasks.

In this work, we propose *Visual Periodic Task Learner* (ViPTL): a method with an explicit periodic policy representation that enables efficient robot learning of periodic robot manipulation skills from a single, visual human demonstration. As a policy representation we adapt rhythmic *Dynamic Movement Primitives* (rDMPs) [11] which model cyclic motion, with shift parameters that account for motion that translates over multiple periods. The benefits of using rhythmic DMPs are two-fold: (1) they succinctly represent

† Work done while the author was an intern at IPRL lab at Stanford.



Fig. 1: Examples of everyday periodic tasks. Our robot learns to imitate wiping, stirring and winding from human video demonstrations (top row).

periodic manipulation policies, enabling efficient learning – the problem of learning a long-horizon periodic task is reduced to learning a single-period motion of the task; (2) the learned rhythmic DMP can be repeated any number of times and the speed and amplitude of the motion can be adjusted at run time by a simple change of a parameter.

Typically, DMPs are trained on robot trajectories demonstrated through kinesthetic teaching or teleoperation. However, this can be non-intuitive to non-expert users. Instead, we use human video demonstrations that do not contain trajectories but are easier to record. To train rDMPs from such videos, we learn a keypoint detector that identifies consistent keypoints across human demonstrations and robot executions. Based on this, we can now evaluate the similarity between a robot execution and the human demonstration - a quantity we aim to maximize. Keypoint models do not assume rigid objects and are therefore well suited for challenging manipulation tasks considered in this paper that involve deformable objects and granular material. We learn the keypoint model from task-agnostic play data and leverage periodicity estimation from [12] to break down the demonstrated task into single-period components. We use Bayesian optimization (BO) to optimize similarity of the robot motion to the human demonstration. To focus BO on promising regions, we create 'imagined' trajectories from segments of robot play data that serve as initial candidates for BO.

We quantitatively evaluate the proposed method in both simulation and on a real robot with three manipulation tasks: table wiping, rope winding, and food stirring (see Figure 1). We show that our approach can successfully learn challenging periodic manipulation tasks that involve deformable and granular objects from a single human demonstration within 50 robot trials. Our comparisons to existing approaches and ablations show how our perception and optimization modules contribute to the overall success of the method.

<sup>1.</sup> Machine Learning Department, Carnegie Mellon University; 2. Department of Mechanical Engineering, Stanford University; 3. Department of Computer Science, Stanford University; 4. Facebook AI Research.

<sup>‡</sup> Supported by the National Science Foundation grant No.2030859 to the Computing Research Association for the CIFellows Project.

Contacts: jingyuny@andrew.cmu.edu, rika.antonova@stanford.edu



Fig. 2: Overview of the proposed approach. Our method is composed of two modules: the unsupervised visual representation module and the parameter optimization module. The unsupervised module learns a model for keypoint correspondences between the motions of the objects in the human demonstration and the robot trials. The parameter optimization module uses active learning to adjust the parameters of the rhythmic DMP controllers.

### **II. RELATED WORK**

### A. Modeling Periodic Motion

Periodic motion is a common component of robotic tasks. Various methods have been proposed to model such motion. Early literature in robotics and neuroscience used limit cycles and central pattern generators to model periodic motions for locomotion [13]–[16]. Recently, pattern generators have also been used with robotic manipulators [17], [18]. However, limit cycle formulations are not easily amenable to learning arbitrary periodic trajectories. In such cases, dynamic movement primitives (DMPs) [11] can provide the needed flexibility, and ease of use with learning-based approaches. DMPs have been used for periodic manipulation tasks, writing and wiping being the most common [19], [20]. In comparison, our proposed work also learns to do winding and stirring tasks, but from visual demonstration that are not annotated with human hand poses. Fourier movement primitives (FMPs) [21] are an extensions of DMPs using Fourier series as basis functions. While our system is agnostic to the specific choice of periodic parameterization of the control policies, here we use rhythmic DMPs. In future work, we will investigate alternative representations such as FMPs.

### B. Periodicity Estimation

There has been a significant interest in estimating periodicity in the computer vision community. Prior works use Fourier analysis [22]–[24], singular value decomposition [25], or peak detection [26] to detect repetition by converting the motion in videos to one-dimensional signals. Recent works propose detecting non-stationary repetitive motion using wavelet transforms [27], 3D convolution networks [28], and self-similarity between video frames [12].

In this work, we use RepNet [12] for periodicity estimation. We find that once trained on the Countix dataset in [12], RepNet can successfully decompose human demonstrations of various manipulation tasks into single-period segments without any further finetuning.

### C. Learning from Human Demonstrations

Several works in learning from human video demonstrations propose using image-to-image translation to transform human demos to robot executions [1]–[4]. However, these require a large amount of training data. Recent works [5], [6] leverage action recognition models, such as the action classifiers trained on the 20BN Something-Something dataset [29], to identify whether the robot is performing the desired task. However, while these classifiers are useful for identifying the class of motions for short interactions, we show that they do not retain enough information to analyze tasks with longer duration and multiple repetitions. Furthermore, the ability of these methods to handle highly deformable objects, such as cloth and ropes, has not been studied yet.

Our approach uses a small amount of task-agnostic, unpaired and unlabeled 'play' data [30] to train a keypoint model that makes it possible to quantitatively compare human demos and robot executions. 'Play' data is useful, because it can be collected without supervision and using a task-agnostic, randomized policy. However, unlike [30] that explores using hours of such data, we focus on a much more data-efficient alternative. We collect 10 minutes of human 'play' data and 10 minutes of robot 'play' data (unpaired), and then use a single human demonstration video to infer the appropriate parameters for the robot control policy.

# D. Sample-efficient Robot Learning

We aim to imitate a visual human demo on a robot with high sample-efficiency. Prior works have explored modelbased methods [31]–[33] and self-supervised exploration algorithms [34]–[39] to improve sample efficiency, but these methods often require much more data on the robot and do not directly generalize to visual imitation learning. Some works utilize large-scale training in simulation and transfer the learned policies to the real robot [40], [41]. Since our tasks include hard-to-simulate deformable and granular objects, sim2real is not applicable in our setup. In contrast, we use Bayesian Optimization (BO) to optimize parameters of a rhythmic DMP. BO is capable of learning the demonstrated manipulation skills within 50 trials on the real robot.

# **III. PRELIMINARIES**

We consider the problem of learning periodic manipulation skills from a single human demonstration. We assume that the human demonstrates a periodic task that is performed for at least 2 periods. Given a human demonstration  $V_H = (I_H^1, \ldots, I_H^{T_H})$  as a sequence  $V_H$  of  $T_H$  RGB image frames



Fig. 3: Unsupervised keypoint learning from play data. The play data consists of unpaired, unlabeled and task-agnostic human and robot motion recorded from the same viewpoint. It is used to train a keypoint model that finds consistent keypoints across human and robot demonstrations. This allows to compute a keypoint-based distance between human and robot videos suitable for guiding the search for the best matching robot motion. Robot trajectories are included in the robot play data, but they are not used for keypoint learning purposes.

 $I_H^t$ , the robot is allowed to execute a total of 50 trials in the environment to learn the demonstrated skill. In each trial, the robotic agent executes a trajectory  $\tau_R = (x_R^1, \ldots, x_R^T)$ where  $x^t$  denotes the robot end-effector position at timestep t. We denote the corresponding execution video of this executed trajectory as  $V_R = (I_R^1, \ldots, I_R^T)$ , where  $I_R^t$  denotes the camera image at timestep t. We measure the similarity between the object motion in the human demonstration and in the video of the robot execution using consistent keypoints across the two videos (see Section IV-A). In this work, we combine BO and DMPs to maximize this similarity score, and present our method in Section IV. Below we provide the technical background for BO and DMPs.

### A. Bayesian Optimization

In Bayesian optimization (BO), an optimization problem is viewed as finding parameters **w** that optimize some objective function  $f_{BO}(\mathbf{w}) : f_{BO}(\mathbf{w}^*) = \max_{\mathbf{w}} f_{BO}(\mathbf{w})$ .  $f_{BO}$  is commonly modeled with a Gaussian process (GP):  $f_{BO}(\mathbf{w}) \sim \mathcal{GP}(m(\mathbf{w}), k(\mathbf{w}_i, \mathbf{w}_j))$ . At each trial, to select the next promising candidate **w**, BO optimizes an acquisition function, e.g. the Upper Confidence Bound (UCB) [42], which explicitly balances exploration (high posterior uncertainty) vs exploitation (high posterior mean estimate):  $UCB(\mathbf{w}) = m(\mathbf{w}) + \beta Var(\mathbf{w})$ . The kernel defines a similarity function on the search space. RBF kernel is a common choice:  $k(\mathbf{w}_i, \mathbf{w}_j) = \sigma_k^2 \exp(-\frac{1}{2} ||\mathbf{w}_i \cdot \mathbf{w}_j||_2^T \operatorname{diag}(\mathbf{l})^{-2} ||\mathbf{w}_i \cdot \mathbf{w}_j||_2)$ , where  $\sigma_k^2$  and I are signal variance and a vector of length scales, respectively. In practice,  $\sigma_k^2$  is a hyperparameter optimized automatically by maximizing marginal likelihood.

# B. Dynamic Movement Primitives (DMPs)

DMPs are trajectory generators whose parameters can be learned from demonstrations of desired robot end-effector trajectories. They combine linear fixed-point attractors with non-linear function approximators to encode complex trajectories, while maintaining convergence guarantees. We refer readers to [43] for a detailed overview.

The *transformation system* of a DMP consists of a damped linear feedback term, and a forcing function f:

$$\tau \dot{z} = \alpha_z \left(\beta_z (g - x) - z\right) + f \; ; \quad \tau \dot{x} = z \; , \tag{1}$$



Fig. 4: Our Bayesian optimization pipeline.

where x is position, g is the goal,  $\alpha_z$  and  $\beta_z$  are constants,  $\tau$  is a temporal scaling factor, z is the scaled velocity, and the output of the transformation system is the scaled acceleration  $\dot{z}$ . The second component of DMPs is a *canonical system* which replaces time, and enables scaling the trajectory to different time lengths. The canonical system is different between rhythmic and discrete DMPs; in the discrete case it represents 'time left' and goes to 0 at the end of the motion, while in the rhythmic case it represents the time from the start, and goes up linearly. Specifically, in rhythmic DMPs, the first-order canonical system  $\tau \dot{\phi} = 1$  encodes the phase  $\phi$ , increasing linearly as motion progresses.

In rhythmic DMPs, the forcing function f is parameterized by  $\phi$  and consists of cyclic basis functions:

$$f = \frac{\sum_{i} \Psi_{i} w_{i}}{\sum_{i} \Psi_{i}} r ; \Psi_{i} = \exp\left(h_{i}\left(\cos\left(\phi - c_{i}\right) - 1\right)\right), \quad (2)$$

where  $\Psi$  is a function of the canonical system, and the weights  $w_i$  are commonly learned using locally weighted regression [44]. The cyclic nature of basis functions ensures that the transformation system yields cyclic motion, as the canonical system unrolls. Typically, the goal g of a rhythmic DMP is set to the mean of the demonstration trajectory and kept fixed. Discrete DMPs are shown to generalize well to changing goals, and [43] present ways to continuously change goals to new locations without causing discontinuity in the acceleration  $\dot{z}$ . We adapt [43] to smoothly move the goals of rhythmic DMPs between executions. This continuously modulates the mean point of the limit cycle of the DMP, allowing us to model motions that are mostly cyclic, but slightly shifting over time, e.g. as in wiping a surface.

# IV. METHOD

Our framework is composed of two parts: (1) a representation learning module, where a keypoint detection model is trained to extract consistent keypoints from independently collected and non-task-specific human and robot play data; (2) a parameter optimization module, where BO searches for a rhythmic DMP that when executed produces a robot video that matches the human demo in terms of the detected keypoints. These modules are detailed in Sections IV-A and IV-B, respectively. Figure 2 shows an overview.

### A. Unsupervised Keypoint Learning from Play Data

To learn a manipulation skill from a human demo, we need a way to evaluate the similarity between the demo and robot execution. To learn such a similarity score, we assume that the agent has access to a small amount of human and robot play data. Play data is a dataset of self-guided, task-agnostic, and diverse interactions. The human play data  $\mathcal{D}_H = (I_{HP}^1, \ldots, I_{HP}^{T_{HP}})$  is a sequence of  $T_{HP}$  unlabeled RGB image frames, while the robot play data is a sequence of  $T_{RP}$  RGB image frames  $\mathcal{D}_R = (I_{RP}^1, \ldots, I_{RP}^{T_{HP}})$  accompanied by robot end-effector positions  $x_{RP}^1, \ldots, x_{RP}^{T_{RP}} \in \mathbb{R}^3$ . Note that  $\mathcal{D}_H$  and  $\mathcal{D}_R$  are unpaired and independent.

To acquire a visual representation for the manipulated object that is invariant to change of agent between human and robot, we adopt a variation of the Transporter architecture [45] – an unsupervised keypoint detection model that learns to generate temporally consistent keypoints  $\Psi^*(I)$  on image input I. The learning process is illustrated in Figure 3. Humans and robots may move their hands very differently to generate the same object movement. To make sure that the keypoint model  $\Psi^*$  allocates keypoints to the manipulated objects and not to the human and robot hand, we mask these areas in the reconstruction loss [45] when training the keypoint model. This is achieved by using the commonly available hand detectors and depth filtering methods, respectively, when computing the reconstruction loss that the Transporter is trained on. With this, the keypoint model is more likely to place keypoints on the objects, and be robust to different visual appearance of human and robot hands.

After the keypoint model is trained, we process the human demo  $V_H$  and robot execution  $V_R$  to produce sub-sampled videos  $V_{H'} = \{I_{H'}^i\}_{i=1}^{N_s}$  and  $V_{R'} = \{I_{R'}^i\}_{i=1}^{N_s}$  that both have length  $N_s$ . We then define the distance between the human demo and the robot execution as:

$$D(V_H, V_R) = \frac{1}{N_s N_k} \sum_{i=1}^{N_s} \|\Psi^*(I_{H'}^i) - \Psi^*(I_{R'}^i)\|_1, \quad (3)$$

where  $\|\cdot\|_1$  denotes  $L_1$  norm and  $\Psi^*(I)$  denotes the locations of the  $N_k$  detected keypoints of an image I normalized to range [0, 1]. Note that by using the above distance function, we are optimizing robot trajectories to align with the given human demo (i.e. when the human is p% into the demo, the robot also aims to be roughly p% into the execution).

### B. Few-shot Motion Optimization with BO

With the learned keypoint representation for comparing human demos and robot executions, the problem of learning a periodic robot manipulation task from human demos can be reduced to searching for trajectories that, when executed on the robot, produce execution videos that have low distance to the provided human demo. This means that we can use the keypoint-based distance between a robot execution and a human demonstration as the objective for our optimization problem formulated in Section III. The remaining problem is thus how to efficiently find trajectories to be executed on the robot that best imitate the human demo. Training robot skills on large amounts of simulated data and then using sim2real techniques to transfer the skill to a real system is a common robot learning paradigm. However, deformable and granular objects are hard to simulate realistically and therefore pose a challenge for sim2real transfer. Thus, we propose a method to directly optimize the motion policy on the real robot.



Fig. 5: Imagined trajectories as initial candidates for BO.

1) Periodicity Estimation with RepNet: To imitate periodic manipulation skills shown in the human demo, we need to first determine the periodicity of this demo. We use RepNet [12] – an approach that can estimate when and how often a periodic task is repeated in a video to estimate periodicity. We observe that the RepNet model trained on the Countix dataset can reliably predict the periodicity of the human demos that we consider. So, we use the trained model (without any finetuning) to predict the number of periods  $n_H^{rep} = \text{RepNet}(V_H)$  of the human demo  $V_H$ .

2) Motion Optimization with BO: We propose to optimize single-period waypoints as BO parameters, then use rhythmic DMPs to fit a smooth trajectory and unroll it for multiple periods, as illustrated in Figure 4. Concretely, BO optimizes single-period waypoints:  $\mathbf{w} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_L]$ . To execute a BO sample, we apply a cubic smoothing to the sampled waypoints, then fit a rhythmic DMP to this smooth trajectory and execute it for  $n_H^{rep}$  periods, with the goal g of the DMP shifting  $\mathbf{v}_L - \mathbf{v}_1$  between two consecutive periods.

During conventional BO, the candidates for the first few trials are sampled at random. In the subsequent trials, an acquisition function samples N candidates at random from the search space, then evaluates their posterior mean and variance to select the most promising next candidate. However, in high-dimensional spaces (above 10D) it is unlikely to sample a well-performing candidate randomly. Even with waypoints as the search space for BO, the space is very large. A leading BO method that recently reported success in high-dimensions [46] was not able to reliably succeed on our tasks within 50 trials. Hence, the need to further improve the data efficiency of BO. Our insight is that robot play data contains meaningful interactions that can help BO to focus on the promising regions of the search space. To effectively use this data, we first generate a set of  $N_s$  play data segments  $S = \{\tau_s^1, \tau_s^2, \dots, \tau_s^{N_s}\}$ , where each element  $\tau_s^i = x_{RP}^{t_i:t_i+T_s}$  is a randomly sampled fixed length trajectory in the robot play data of length  $T_s$ . Then, we generate 'imagined trajectories' by rejection sampling segment sequences  $\tau_I = (\tau_s^{k_1}, \dots, \tau_s^{k_m})$  such that the end position  $\tau_s^{k_i, T_s}$  of each segment  $\tau_s^{k_i}$  is less than  $d_{\text{seg}}$  away from the start position  $\tau_s^{k_{i+1}, 1}$  of the next segment  $\tau_s^{k_{i+1}}$ in L2 distance. Then, we can find the corresponding image frames of  $\tau_I$  to construct an 'imagined' video  $V_I$ , and this video can be evaluated using the objective score function  $D(V_H^*, V_I)$ , where  $V_H^*$  is a single-period demo trimmed from the original human demo according to RepNet period split. We then select the top  $N_{\text{imagined}}$  trajectories with the highest



Fig. 6: The left part shows the objectives used in the methods we compare, then explains our evaluation metric – 'performance' in plots on the right. Plots (a)-(c) show the performance of the competing methods in the 3 tasks we consider, in simulation. For the *Direct Imitation* baseline, executions are fixed, no finetuning between trials. We include two versions of this baseline: (1x) – trained from the robot play data; (2x) – trained from twice the amount of that data, so that the total size of training data exceeds the size of robot play data + the 50 trials of interactions. The performance saturates, showing no benefit from additional training data (red lines match). For the other 3 methods, we execute 50 trials for each run. For every method, we do 3 runs using 3 random seeds. The solid lines denote the mean of performance across 3 seeds; the shaded areas denote the standard deviation of the performance values. The MBIL line denotes model-based imitation learning baseline, which learns dynamics using keypoints as states and uses MPC for planning.

estimated scores and construct a set of *initial candidates*:  $\{\mathbf{w}_i\}_{i=1...N_{\text{imagined}}}$ . Each candidate  $\mathbf{w}_i$  is represented by a set of waypoints sub-sampled from the imagined trajectory. We illustrate this process in Figure 5.

To warm-start BO, we sample the candidates for the first few trials from  $\{\mathbf{w}_i\}_{i=1...N_{\text{imagined}}}$ . By construction, these represent the waypoints of the trajectories that have a high alignment with the human demo for the 1st period. In the subsequent trials, we augment the pool of the candidates considered by the acquisition function by sampling in the regions close to these initial candidates. With that, the acquisition function can help us focus the search on the regions close to the initial candidates, but is not restricted to these regions. Hence, we avoid placing hard restrictions on the search space of BO based on prior information. As a result, our BO extension retains theoretical guarantees of BO, such as consistency and regret bounds.

# V. EXPERIMENTS

In our experiments, we aim to answer the following questions: (1) does our framework successfully learn to perform periodic tasks from a single human demonstration; (2) does our proposed framework perform better than methods that do not exploit periodicity in the target task; (3) is our proposed keypoint-based representation more suitable for learning from human demonstrations than other representations (e.g. latent vectors generated by 20BN classifiers)?

We consider 3 challenging manipulation tasks: (1) Table Wiping, where the objective is to wipe a rectangular table surface with a cloth using back-and-forth motions, shifting to cover all the visible the area of the table; (2) Rope Winding, where the objective is to wrap a rope around a fixed spool by repeating circular winding motion several times; (3) Food Stirring, where the objective is to stir granular objects in a tray with a spatula/spoon.

*Metric:* In our method, the baselines and ablations are each optimizing a different objective function. Therefore, we define a performance metric that is comparable across the different approaches. We prepare an exemplary robot trajectory  $\tau_E = (x_E^1, \ldots, x_E^{T_E})$  that, when executed, produces the same effect on the objects being manipulated as

the human demo. During execution, the robot will execute a trajectory  $\tau_R = (x_R^1, \dots, x_R^T)$ , where  $x^t$  denotes the robot end-effector position at timestep t. Performance of an execution that produces robot trajectory  $\tau_R$  is evaluated by the similarity between  $\tau_E$  and  $\tau_R$ . More concretely, the performance is computed to be  $\kappa(-\|\tau_E - \tau_{R'}\|_1)$ , where  $\tau_{R'}$  is a sub-sampled trajectory of  $\tau_R$  with length  $T_E$  and  $\kappa$  is a linear transformation that ensures the score is in range [0.0, 1.0]. Note that  $\tau_E$  are only used for evaluation purposes and are not visible to any of the methods. Figures 6and 7 plot this metric for our simulation and hardware experiments. The methods that use BO each optimize a different objective function e.g. keypoint-based objective for *ViPTL (our method)*, cosine distance between latent features for Twentybn Classifier. We include plots of these in the supplementary video to illustrate BO progress over trials.

### A. Simulation Experiments

1) Baselines and ablations: To pick appropriate baselines and ablations, we need methods that can imitate a single visual human demonstration on the robot. Standard imagebased model-free and model-based RL methods [31], [32], [47], [48] cannot operate in this setup because images from human demo and robot execution are visually different. Thus, we use two baselines that use our keypoint-based visual representation as state and an ablation that does not use this visual representation to test whether our visual representation contributes to the final performance of our method.

First we have *Direct Imitation*, which learns a function that maps keypoints at the current timestep to desired robot end-effector positions. This function is trained on robot play data, which contains both keypoints and robot trajectories. To imitate a human demonstration, we use keypoints from the demo video frames as input and output desired robot end-effector positions. The resulting robot trajectory is then executed by fitting a DMP to the predicted robot positions. This baseline studies the use of BO versus a trained neural network for optimizing DMP controllers.

Second we have *Twentybn Classifier*, an ablation in which the BO objective is based on the video activity classifier [49]



Fig. 7: Results for the real robot experiments. Plots (a)-(c) show performance comparisons (for details of the methods see Section V-A.1, for evaluation metric details see Figure 6). The right part shows skills learned with our method executing at different scales and numbers of repetitions without retraining.

trained on the 20bn dataset [29]. Both the human demonstration and robot execution video are input to the classifier to obtain two feature vectors from the last hidden layer. The objective function is the cosine distance between these two features. We use this baseline to test if the keypoint-based visual representation is a suitable visual representation for the learning from human demo setup compared to alternatives.

Third we have *MBIL*, a model-based imitation learning baseline that relies on a learned dynamics model of the keypoints. The model is trained on robot play data and updated every episode as new interaction data is collected. This model is then used to plan actions to imitate the human demo. At every timestep, we run single-step model-predictive control (MPC) by sampling 5,000 random actions and executing the action that leads to the smallest keypoint distance to the corresponding human demo frame. This baseline tests if our method outperforms model-based RL methods like [50] that do not model periodicity of the task.

2) Experimental Setup: For all our experiments, we use an image size of  $512 \times 512$ . In the distance function, the number of sub-sampled frames  $N_s$  is set to  $10 \cdot n_H^{rep}$ . To create masks for the human hand, we use the MediaPipe library [51] to detect the hand skeleton from an image frame and use the colors at the joints of the skeleton to construct a color range mask for the hand. We mask out the robot based on the depth readings (since robot pose is known). In BO, we optimize L=7 trajectory waypoints in the wiping and winding, L=5 in stirring. We use UCB [52] with  $\beta = 0.1$  in the acquisition function of BO. We use automatic hyperparameter optimization to find the appropriate length scales of the RBF kernel. When constructing imagined trajectories, we use 10 play data segments of length  $T_s = 10$  each and use a distance threshold of  $d_{\text{seg}} = \frac{1}{6} \lambda_{\text{disp}}$ . We generate 5,000 imagined trajectories and select the top  $N_{\text{imagined}} = 100$  trajectories as initial candidates for BO, and use 10 of these in the first 10 BO trials.

3) Quantitative Results: To evaluate the performance of our framework in comparison to competing methods, we run all methods and baselines for 50 trials in all tasks using 3 different random seeds. The performance of all the methods during BO trials is shown in Figure 6. The *Direct Imitation* and *MBIL* baselines cannot imitate the human demo well in Table Wiping and Rope Winding, since modeling or capturing dynamics of the deformable objects in these tasks is difficult. The suboptimal performance of these two baselines shows that our method outperforms methods that rely on single-step predictions or learning accurate dynamics models and do not exploit periodicity in the target task. The *Twentybn Classifier* baseline achieves limited performance in all three tasks as it lacks precision in the classifier-based distance metric used to compare the given human demo with robot executions. This shows that the keypoint-based cost function is a crucial component of our method that is more suitable for the learning from human demos setup. In contrast, our method (*ViPTL*) is able to achieve good performance in all three tasks, outperforming both the *Direct Imitation* baseline and the *Twentybn Classifier* ablation. We include qualitative results in the supplementary video.

### B. Real Robot Experiments

1) Hardware Setup: Our hardware setup (Figure 7) includes a Kinova Gen3 robot arm with a Robotiq 2F-85 gripper and an Intel RealSense D435 camera. The table workspace measures  $50 \times 43$  cm, and the camera is mounted at the height of 68 cm at one side of the workspace. The camera provides the RGB image data during experiments and is positioned to view the table surface. We use velocity control in the Cartesian (end-effector) space to execute the desired trajectories on the robot.

2) Quantitative Results: We select the Twentybn Classifier ablation and the baseline with the most consistent performance across tasks in simulation experiments (Direct Imitation) to compare with our method (see Figure 7). Our method is able to quickly find high-scoring points due to an effective warm-start and fine-tune the generated motion, leading to consistently improving performance through out the 50 trials, while the baselines are unable to catch up with the performance of our method for similar reasons as mentioned in simulated experiments.

### VI. CONCLUSION

We introduced Visual Periodic Task Learner (ViPTL), a method for representing periodic manipulation policies and efficiently learning them from a single human demonstration. We show that ViPTL succeeds on three robot manipulation tasks that involve deformable and granular objects. This work opens the opportunity to benefit from the periodic structure of many tasks commonly seen in everyday life. In future work, we intend to extend our method to handle the initial stages of the tasks, such as grasping and other transient motions.

### REFERENCES

- Y. Liu, A. Gupta, P. Abbeel, and S. Levine, "Imitation from observation: Learning to imitate behaviors from raw video via context translation," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1118–1125.
- [2] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine, "Avid: Learning multi-stage tasks via pixel-level translation of human videos," *arXiv preprint arXiv:1912.04443*, 2019.
- [3] P. Sharma, D. Pathak, and A. Gupta, "Third-person visual imitation learning via decoupled hierarchical controller," *arXiv preprint* arXiv:1911.09676, 2019.
- [4] H. Xiong, Q. Li, Y.-C. Chen, H. Bharadhwaj, S. Sinha, and A. Garg, "Learning by watching: Physical imitation of manipulation skills from human videos," arXiv preprint arXiv:2101.07241, 2021.
- [5] L. Shao, T. Migimatsu, Q. Zhang, K. Yang, and J. Bohg, "Concept2robot: Learning manipulation concepts from instructions and human demonstrations," 2020.
- [6] A. S. Chen, S. Nair, and C. Finn, "Learning generalizable robotic reward functions from" in-the-wild" human videos," arXiv preprint arXiv:2103.16817, 2021.
- [7] D. Xu, S. Nair, Y. Zhu, J. Gao, A. Garg, L. Fei-Fei, and S. Savarese, "Neural task programming: Learning to generalize across hierarchical tasks," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 3795–3802.
- [8] T. Silver, K. R. Allen, A. K. Lew, L. Kaelbling, and J. Tenenbaum, "Few-shot bayesian imitation learning with logic over programs," *ArXiv*, vol. abs/1904.06317, 2019.
- [9] S. Nair and C. Finn, "Hierarchical foresight: Self-supervised learning of long-horizon tasks via visual subgoal generation," *ArXiv*, vol. abs/1909.05829, 2020.
- [10] A. Hundt, B. Killeen, H. Kwon, C. Paxton, and G. Hager, "good robot!": Efficient reinforcement learning for multi-step visual tasks via reward shaping," *ArXiv*, vol. abs/1909.11730, 2019.
- [11] A. Ijspeert, J. Nakanishi, and S. Schaal, "Learning rhythmic movements by demonstration using nonlinear oscillators," *IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 1, pp. 958–963, 2002.
- [12] D. Dwibedi, Y. Aytar, J. Tompson, P. Sermanet, and A. Zisserman, "Counting out time: Class agnostic video repetition counting in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10387–10396.
- [13] L. Jalics, H. Hemami, and Y.-F. Zheng, "Pattern generation using coupled oscillators for robotic and biorobotic adaptive periodic movement," in *Proceedings of International Conference on Robotics and Automation*, vol. 1. IEEE, 1997, pp. 179–184.
- [14] S. Kajita, F. Kanehiro, K. Kaneko, K. Fujiwara, K. Yokoi, and H. Hirukawa, "A realtime pattern generator for biped walking," in *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, vol. 1. IEEE, 2002, pp. 31–37.
- [15] J. F. Yang, T. Lam, M. Y. Pang, E. Lamont, K. Musselman, and E. Seinen, "Infant stepping: a window to the behaviour of the human pattern generator for walking," *Canadian journal of physiology and pharmacology*, vol. 82, no. 8-9, pp. 662–674, 2004.
- [16] L. Righetti and A. J. Ijspeert, "Pattern generators with sensory feedback for the control of quadruped locomotion," in 2008 IEEE International Conference on Robotics and Automation, 2008, pp. 819– 824.
- [17] M. Thor and P. Manoonpong, "A fast online frequency adaptation mechanism for cpg-based robot motion control," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3324–3331, 2019.
- [18] P. Oikonomou, M. Khamassi, and C. S. Tzafestas, "Periodic movement learning in a soft-robotic arm," in 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 4586–4592.
- [19] J. Ernesti, L. Righetti, M. Do, T. Asfour, and S. Schaal, "Encoding of periodic and their transient motions by a single dynamic movement primitive," in 2012 12th IEEE-RAS International Conference on Humanoid Robots (Humanoids 2012), 2012, pp. 57–64.
- [20] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal, "Dynamical Movement Primitives: Learning Attractor Models for Motor Behaviors," *Neural Computation*, vol. 25, no. 2, pp. 328–373, 02 2013.
- [21] T. Kulak, J. Silvério, and S. Calinon, "Fourier movement primitives: an approach for learning rhythmic robot skills from demonstrations," in *Proc. Robotics: Science and Systems (RSS)*, 2020.

- [22] O. Azy and N. Ahuja, "Segmentation of periodically moving objects," in 2008 19th International Conference on Pattern Recognition. IEEE, 2008, pp. 1–4.
- [23] R. Cutler and L. S. Davis, "Robust real-time periodic motion detection, analysis, and applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 781–796, 2000.
- [24] E. Pogalin, A. W. Smeulders, and A. H. Thean, "Visual quasiperiodicity," in 2008 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2008, pp. 1–8.
- [25] D. Chetverikov and S. Fazekas, "On motion periodicity of dynamic textures." in *BMVC*, vol. 1. Citeseer, 2006, pp. 167–176.
- [26] A. Thangali and S. Sclaroff, "Periodic motion detection and estimation via space-time sampling," in 2005 Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION'05)-Volume 1, vol. 2. IEEE, 2005, pp. 176–182.
- [27] T. F. Runia, C. G. Snoek, and A. W. Smeulders, "Real-world repetition estimation by div, grad and curl," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9009– 9017.
- [28] H. Zhang, X. Xu, G. Han, and S. He, "Context-aware and scaleinsensitive temporal repetition counting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 670–678.
- [29] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, *et al.*, "The" something something" video database for learning and evaluating visual common sense," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5842–5850.
- [30] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet, "Learning latent plans from play," in *Conference on Robot Learning*. PMLR, 2020, pp. 1113–1132.
- [31] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba, "Mastering atari with discrete world models," arXiv preprint arXiv:2010.02193, 2020.
- [32] F. Ebert, C. Finn, S. Dasari, A. Xie, A. Lee, and S. Levine, "Visual foresight: Model-based deep reinforcement learning for vision-based robotic control," arXiv preprint arXiv:1812.00568, 2018.
- [33] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak, "Planning to explore via self-supervised world models," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8583–8592.
- [34] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *International conference* on machine learning. PMLR, 2017, pp. 2778–2787.
- [35] J. Schmidhuber, "A possibility for implementing curiosity and boredom in model-building neural controllers," in *Proc. of the international conference on simulation of adaptive behavior: From animals to animats*, 1991, pp. 222–227.
- [36] A. S. Klyubin, D. Polani, and C. L. Nehaniv, "Empowerment: A universal agent-centric measure of control," in 2005 ieee congress on evolutionary computation, vol. 1. IEEE, 2005, pp. 128–135.
- [37] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," arXiv preprint arXiv:1802.06070, 2018.
- [38] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos, "Unifying count-based exploration and intrinsic motivation," *Advances in neural information processing systems*, vol. 29, pp. 1471–1479, 2016.
- [39] D. Yarats, R. Fergus, A. Lazaric, and L. Pinto, "Reinforcement learning with prototypical representations," arXiv preprint arXiv:2102.11271, 2021.
- [40] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, 2017, pp. 23–30.
- [41] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, "Sim-to-real via simto-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 627–12 637.
- [42] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger, "Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design," arXiv preprint arXiv:0912.3995, 2009.
- [43] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal,

"Dynamical movement primitives: learning attractor models for motor behaviors," *Neural computation*, vol. 25, no. 2, pp. 328–373, 2013.

- [44] S. Schaal and C. G. Atkeson, "Constructive incremental learning from only local information," *Neural computation*, vol. 10, no. 8, pp. 2047– 2084, 1998.
- [45] T. Kulkarni, A. Gupta, C. Ionescu, S. Borgeaud, M. Reynolds, A. Zisserman, and V. Mnih, "Unsupervised learning of object keypoints for perception and control," arXiv preprint arXiv:1906.11883, 2019.
- [46] D. Eriksson, M. Pearce, J. Gardner, R. D. Turner, and M. Poloczek, "Scalable global optimization via local Bayesian optimization," in *Advances in Neural Information Processing Systems*, 2019, pp. 5496– 5507.
- [47] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, *et al.*, "Qtopt: Scalable deep reinforcement learning for vision-based robotic manipulation," *arXiv preprint arXiv:1806.10293*, 2018.
- [48] Y. Chebotar, K. Hausman, Y. Lu, T. Xiao, D. Kalashnikov, J. Varley, A. Irpan, B. Eysenbach, R. Julian, C. Finn, *et al.*, "Actionable models: Unsupervised offline reinforcement learning of robotic skills," *arXiv* preprint arXiv:2104.07749, 2021.
- [49] L. Wang, Z. Tong, B. Ji, and G. Wu, "Tdn: Temporal difference networks for efficient action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1895–1904.
- [50] L. Manuelli, Y. Li, P. Florence, and R. Tedrake, "Keypoints into the future: Self-supervised correspondence in model-based reinforcement learning," arXiv preprint arXiv:2009.05085, 2020.
- [51] C. Lugaresi, J. Tang, H. Nash, C. McClanahan, E. Uboweja, M. Hays, F. Zhang, C.-L. Chang, M. G. Yong, J. Lee, *et al.*, "Mediapipe: A framework for building perception pipelines," *arXiv preprint arXiv:1906.08172*, 2019.
- [52] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger, "Gaussian process optimization in the bandit setting: No regret and experimental design," *arXiv preprint arXiv:0912.3995*, 2009.