

# K-VIL: Keypoints-based Visual Imitation Learning

Jianfeng Gao, Zhi Tao, Noémie Jaquier, and Tamim Asfour

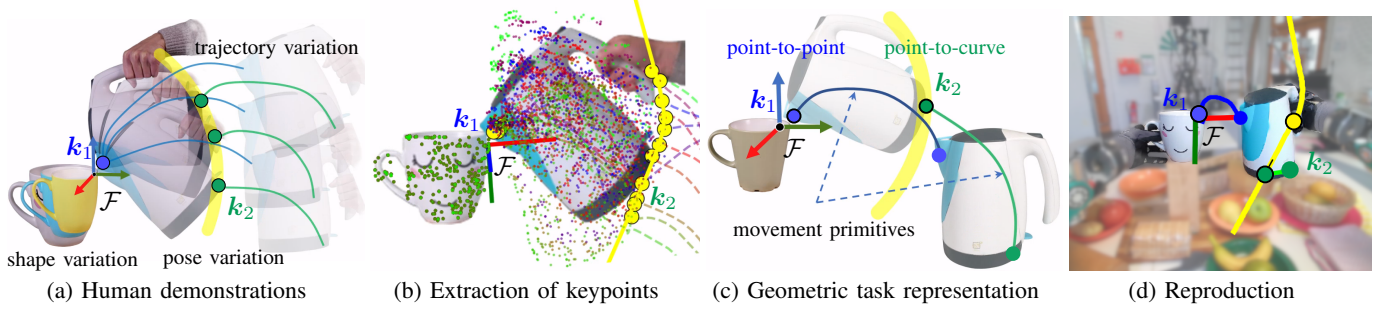


Fig. 1: Overview of the K-VIL approach. (a) Human demonstration videos of manipulation actions involving categorical objects with shape, pose, and trajectory variations. (b) Sampling of dense candidate points from the object surface. (c) Extraction of *sparse keypoints*  $k_1, k_2$  subject to certain types of *geometric constraints* (point-to-point and point-to-curve), their associated *local frames*  $\mathcal{F}$  and the *movement primitives* which represent the demonstrated keypoint motions. (d) Adaptation of the learned generalizable geometric task representation to a new scene, and execution by the robot.

**Abstract**—Visual imitation learning provides efficient and intuitive solutions for robotic systems to acquire novel manipulation skills. However, simultaneously learning geometric task constraints and control policies from visual inputs alone remains a challenging problem. In this paper, we propose the *keypoint-based visual imitation learning* (K-VIL) approach that automatically extracts sparse, object-centric, and embodiment-independent task representations from a small number of human demonstration videos. The task representation is composed of keypoint-based geometric constraints on principal manifolds, their associated local frames, and the movement primitives that are then needed for the task execution. Our approach is capable of extracting such task representations from a single demonstration video, and of incrementally updating them when new demonstrations are available. To reproduce manipulation skills using the learned set of prioritized geometric constraints in novel scenes, we introduce a novel keypoint-based admittance controller. We evaluate our approach in several real-world applications, showcasing its ability to deal with cluttered scenes, viewpoint mismatch, new instances of categorical objects, and large object pose and shape variations. Our evaluation demonstrates the efficiency and robustness of our approach in both one-shot and few-shot imitation learning settings. Videos and source code are available at <https://sites.google.com/view/k-vil>.

**Index Terms**—Learning from Demonstration; Visual Learning; Manipulation Planning; Learning of Geometric Constraints.

## I. INTRODUCTION

This work has been supported by the German Federal Ministry of Education and Research (BMBF) under the project OML and by the Carl Zeiss Foundation under the project JuBot. (Corresponding author: Jianfeng Gao)

The authors are with the Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Karlsruhe, Germany. E-mails: {jianfeng.gao, noemie.jaquier, asfour}@kit.edu; zhitao.robotics@gmail.com.

Published version: <https://ieeexplore.ieee.org/document/10189175>

**O**BSERVATIONAL learning, i.e., the ability to develop new skills from observed actions and their outcome, is an important learning mechanism in our daily lives [1]–[3]. For example, by watching a few videos showing people pouring water from a kettle into different teacups (as in Fig. 1a), we can easily learn “what” a pouring task is and “how” to perform it. From a computational point of view, the spout and bottom of the kettle can be represented by two keypoints  $k_1, k_2$ . As shown in Fig. 1c, a pouring task then simply consists in aligning the spout  $k_1$  with a point above the rim of the cup (point-to-point constraint) and similarly aligning the bottom  $k_2$  with a curve that controls the kettle’s angle of inclination (point-to-curve constraint). Such sets of keypoints and keypoint-based geometric constraints can generally be used to represent daily manipulation tasks, i.e., to parameterize the motion of their functional parts relative to some local frames of reference. Moreover, these keypoints and local frames can also be associated with local visual features of object functional parts.

In this paper, we propose to exploit such task representations to teach manipulation skills to robots from video demonstrations. We additionally aim for generalizable skills, which can be reused in novel scenes (see Fig. 1d for an example). In this context, three main challenges arise, namely, (i) the detection and efficient extraction of task-relevant keypoints on objects; (ii) the definition of generalizable and embodiment-independent task representations; and (iii) the reproduction of the demonstrated task and its adaptation to new scenes. To address the first two challenges, we first *densely* sample a set of candidate points from the object mask provided by Mask R-CNN [4] and leverage the correspondence detection of Dense Object Net (DON) [5] to

track their motions in the demonstration videos and obtain their 3D positions (see Section III for a short background). Then, we jointly extract a set of *sparse keypoints* and a set of *keypoint-based geometric constraints* representing the task, as shown in Fig. 1b (see Section IV). To do so, we exploit principal manifold estimation algorithms (PME) [6], which are intrinsically more data- and time-efficient than approaches based on supervised learning [7], [8] and reinforcement learning (RL) [9], [10]. The resulting geometric constraints are expressed relative to local frames defined on target objects and are easily adjustable to pose and shape variations of the target object. As shown in previous works [7], [9], [11], such object-centric task representations facilitate the transfer of manipulation skills between demonstrators and imitators. It also allows our approach to deal with demonstrations provided from different viewpoints. The representation of the task in the form of keypoints and their constraints also enables the use of simpler control policies [12], such as movement primitives, for executing the given task. In this paper, we exploit this property to address our third challenge. Namely, we encode the keypoint motions relative to the corresponding local frames as *via-point movement primitives (VMPs)* [13], which are flexible in terms of temporal scaling and trajectory adaptation while maintaining the demonstrated motion styles. The learned keypoint motions can then be executed on a robot by leveraging our novel keypoint-based admittance controller (Section V). We validate our approach by learning various real-world daily tasks from video demonstrations and reproducing them with a humanoid robot (Section VI). The results show that K-VIL efficiently extracts generalizable manipulation skills, handles viewpoint mismatch, and deals with large pose and shape variations of categorical objects in cluttered scenes.

Our contributions are threefold: (i) We introduce the *Keypoint-based Visual Imitation Learning (K-VIL)* approach for automatic and incremental extraction of sparse, object-centric, viewpoint-invariant, and embodiment-independent task representations. K-VIL extracts task representations from a single demonstration video and improves them as new demonstrations are available. The task representations consist of keypoint-based geometric constraints on principal manifolds, their associated local frames, and the movement primitives required to reproduce the task. (ii) We formulate and learn a large variety of geometric constraints, which allow the proposed task representation to be flexible and efficient. (iii) We propose a novel keypoint-based admittance controller that handles a set of prioritized geometric constraints and allows successful reproductions of the learned task in novel scenes.

## II. RELATED WORK

Visual Imitation Learning (VIL) is a class of imitation learning (IL) frameworks in which only visual sensory input is presented to the imitator. The main challenges of VIL are (i) the detection of visual correspondences between the demonstrator's and imitator's context, i.e., context translation [14], [15]; (ii) the fine-grained understanding of scene structures [16], along with the design of generalizable task representation; and (iii) the design of sample efficient and

scalable control policies. This latter challenge is often tackled along with the former ones, as described next.

### A. Context Translation

Context translation has typically been addressed by training context translators in the demonstrator context to predict the observations in the imitator (e.g., the robot) context. Pixel-level translators were used in [15], [17], [18], [19] to further train RL policies by maximizing the similarity between predicted and received robot observations. Despite the performance of such models, their training is computationally expensive and time-consuming. To improve learning efficiency, Sharma et al. [14] combined a goal-level translator with a task-agnostic control policy, which was trained independently and shared among different tasks. In contrast to these works, K-VIL represents the context via a set of object-centric keypoints and their respective geometric constraints, thus facilitating the context translation between demonstrators and imitators. In addition, by leveraging Mask R-CNN and DON models — which are trained beforehand in a task-agnostic manner and shared among tasks — K-VIL's representations can be acquired from a single or few demonstrations.

### B. Fine-grained Understanding of Scene Structure

The above approaches do not scale to categorical objects as they do not explicitly extract the scene structures with respect to objects and their functional parts. In the literature, the understanding of fine-grained scene structures is mainly achieved through 1) the viewpoint-invariant representation of fine-grained scene features; 2) the extraction of a hierarchy of the scene structure; and 3) the definition of task constraints.

1) *Viewpoint-invariant representation: Dense visual descriptors* such as DON [5] and Neural Descriptor Fields (NDFs) [20] represent fine-grained scene features by detecting dense correspondences of categorical objects, thus allowing point-based representation of object functional parts. However, in [20], [21] [22], access to the robot state space was required in addition to the visual demonstrations, thus violating the purpose of visual imitation. Yang et al. [23] proposed a transporter-based representation learning model to extract keypoints from the task-agnostic human and robot play data. Building the similarity function of such a model requires robot execution videos with a similar view setup as the demonstration videos. The same requirement applies to the approaches presented in [24]–[26] and prevents robots from learning from human demonstrations taken from a very different viewpoint. Sermanet et al. [27] proposed Time-Contrastive Networks (TCN) to learn viewpoint-invariant latent representations of the scene. This approach requires a large number of demonstration videos and robot play videos to build the correspondence between human and robot arms, which makes the approach embodiment-dependent. Similar to our paper, Karnan et al. [28] proposed to leverage task-agnostic keypoint detection algorithms for vehicle navigation tasks. This approach requires storing the demonstration video and searches the closest demonstration image for reward construction. This reduces the number of demonstrations compared to the pixel-level context

translations of [15], [17]–[19]. However, by overlooking the different types of geometric constraints that the keypoints are subject to, this approach suffers from averaging problem similar to [20] (see Section II-B3 for details). In contrast, K-VIL uses dense point-based object representation and correspondence detection to align demonstrations recorded from different viewpoints. This significantly reduces the required number of demonstrations. Moreover, K-VIL explicitly extracts viewpoint- and embodiment-independent scene structure and task constraints, thus addressing the average problem and achieving better extrapolation capability in fine-grained manipulation tasks.

2) *Hierarchy of scene structure.*: The variance across demonstrations was used to efficiently select appropriate local frames from some candidates in several imitation learning frameworks as a solution to extract hierarchical scene structure [29], [30]. Representing the learned task in such local frames was shown to facilitate the transfer of skills between different embodiments and the design of control policies. However, in the absence of visual sensory input, the candidates were manually defined at object level in [29], [30]. In our work, we instead show that combining dense visual descriptors and a variance-based criterion allows for the efficient extraction of keypoints and local frames at a fine-grained level.

3) *Task constraints*: Early works on visual servoing [31]–[33] hand-crafted task constraints as simple geometric constraints (e.g., point-to-point, point-to-line). To represent more complex constraints, Sieb et al. [16] proposed visual entity graphs (VEGs) based on DON to disentangle the scene structure into multiple levels, including objects, parts, and points. A path integral policy was then trained on the similarity loss between the VEGs learned from the demonstrator and the VEGs observed by the robot. The task constraints are implicitly learned in VEGs, similarly to the neural pose descriptors in [20], and are therefore averaged when large shape or pose variations occur in the demonstrations. To address this issue, Jin et al. [12], [34] introduced an explicit representation of geometric constraints using visual geometric skill kernels and graph neural networks, which generalized better to categorical objects. However, this approach requires  $\sim 30$  demonstrations to learn a generalizable representation, since both task correspondences and geometric constraints need to be learned in the graph structure. In this paper, we exploit the correspondence detection of DON and the variation information to jointly extract explicit, sparse keypoints and endow them with geometric constraints of various types (see Section IV-B). This allows us to learn generalizable skills from only a few demonstration videos while alleviating the averaging problem of [16]. Moreover, our task representation allows us to replace the RL policy used, e.g., in [16], with simpler movement primitives that reproduce the demonstrated keypoint motions and adapt to new goal configurations.

### III. BACKGROUND

In this section, we introduce the dense visual correspondence models, the principal manifold estimation algorithm, and the VMPs, which are essential building blocks of K-VIL.

#### A. Dense Visual Correspondence

Dense Object Net (DON) [5] maps an RGB image  $\mathbf{A} \in \mathbb{R}^{W \times H \times C}$  to a dense descriptor image  $\mathbf{A}_{\bar{D}} \in \mathbb{R}^{W \times H \times \bar{D}}$ , where  $W, H, C$  denote the width, height and the number of channels of the image, and  $\bar{D}$  is the dimension of the descriptor space. Therefore, each pixel of the input image is represented by a  $\bar{D}$ -dimensional descriptor  $\mathbf{d} \in \mathbb{R}^{\bar{D}}$ . To train a DON model on an object category, multiple views of posed RGB images of multiple instances of this object category are first collected. The object meshes are then reconstructed using any state-of-the-art scene reconstruction method. The reconstructed meshes are then exploited to automatically acquire object masks and retrieve dense correspondence signals, which are used to train the model. A fully trained DON model maps similar local patches of two images of the categorical objects to patches in the descriptor space with similar descriptors. In other words, the dense visual correspondence between two pixels is detected if the distance between their descriptors is smaller than a certain threshold. For example, the spout of the kettle in different image frames is mapped to similar descriptors.

In this paper, we aim at retrieving dense correspondences among different instances of the same object category. Moreover, once a set of sparse keypoints is extracted, we aim at identifying these keypoints during the task reproduction on new instances of the same object category using their descriptors. To do so, we construct a correspondence function  $\mathbf{p} = f_c(\mathbf{A}, \mathbf{d})$  using the DON model and the camera intrinsic and extrinsic parameters, which can be used to extract the 3D position of a keypoint represented by the descriptor  $\mathbf{d}$ . Similarly to [5], [35], we use 34-layer, stride-8 ResNet as the DON model and set  $\bar{D} = 3$ . We also train Mask R-CNN models [4] with the automatically generated object mask dataset similar to [16]. We refer the interested readers to [5], [35] for additional details of DON models. For the case where the human hand is involved in the demonstrated tasks, we treat it as a special object and utilize a hand keypoint detection algorithm, e.g., MediaPipe [36], which provides more robust correspondence detection on human hands.

#### B. Principal Manifold Estimation (PME)

Within K-VIL, we are interested not only in extracting a set of keypoints, but also in learning the constraints that they satisfy in order to represent the task. Specifically, we use geometric constraints, which allow us to restrict the keypoint target positions, e.g., as in Fig. 1c. In a 3-D space, a simple geometric constraint can be viewed as a low-dimensional manifold corresponding to a point, a line, a plane, a curve, or a surface. In this paper, we leverage the principal manifold estimation (PME) algorithm [6] to uncover the geometric constraints as low-dimensional embedding from a set of 3D points varying in time (obtained via DON). In PME, the principal manifold is defined as a minimum of the functional with a regularity penalty term derived on a Sobolev space. Specifically, the PME algorithm minimizes the loss

$$\mathcal{L}(f, \pi_d) = \mathbb{E} \|\mathbf{x} - f(\pi_d(\mathbf{x}))\|^2 + \lambda \|\kappa_f\|^2, \quad (1)$$

where  $\pi_d : \mathbb{R}^D \rightarrow \mathbb{R}^d$  is the projection index that maps a random  $D$ -dimensional vector  $\mathbf{x}$  onto a  $d$ -dimensional principal

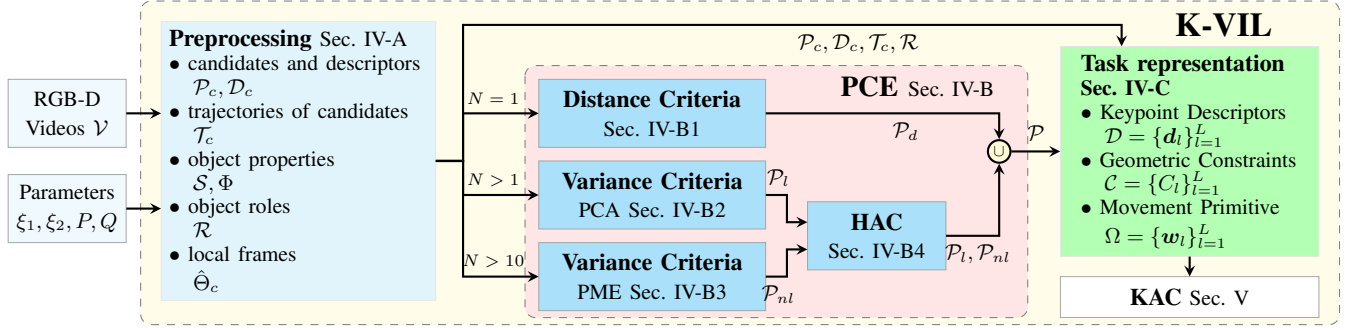


Fig. 2: Overview of K-VIL’s architecture. After preprocessing the demonstration videos, K-VIL jointly extracts a set of sparse keypoints, a set of keypoint-based geometric constraints, and a set of movement primitive parameters that fully represent the task. The robot then leverages the proposed keypoint-based admittance controller to reproduce the task in novel scenes.

manifold with  $d < D$ ,  $f: \mathbb{R}^d \rightarrow \mathbb{R}^D$  is the reconstruction function,  $\|\kappa_f\|^2$  represents the high-dimensional generalization of the total squared curvature of the principal manifold, and  $\lambda \in [0, \infty)$  controls the model complexity. Therefore, the former term of the loss represents the reconstruction error, while the latter regularizes the model to avoid overfitting. Note that PME reduces to linear principal component analysis (PCA) when  $\lambda \rightarrow \infty$ . The linearity and the dimension  $d$  of the principal manifold determine the subspace type. For example, a nonlinear principal manifold of dimension  $d = 1$  is a principal curve and corresponds to a curve constraint. We refer the reader to [6] for the details of the PME algorithm.

### C. Via-point Movement Primitive (VMP)

In addition to extracting the keypoint constraints, we are interested in learning their motions from human demonstration videos. In imitation learning, motions are often represented by movement primitives. Here, we use via-point movement primitives (VMPs) [13]. A VMP combines a linear elementary trajectory  $h_{\text{vmp}}$  with a nonlinear shape modulation  $f_{\text{vmp}}$ , so that

$$y(x) = h_{\text{vmp}}(x) + f_{\text{vmp}}(x) = g + x(y_0 - g) + \psi(x)^\top w,$$

where  $x$  is the canonical variable decreasing linearly from 1 to 0,  $y$  and  $y_0$  are current and start positions, and  $g$  is the target position. The shape modulation term is defined as a linear regression model based on  $N_k$  squared exponential (SE) kernels  $\psi_i(x) = \exp(-h_i(x - c_i)^2)$ ,  $i \in [1, N_k]$ , where  $h_i, c_i$  are pre-defined constants. Similarly to probabilistic movement primitives (ProMP) [37], VMPs assume that the weight parameter  $w \sim \mathcal{N}(\mu_w, \Sigma_w)$  follows a Gaussian distribution, and thus can be learned via maximum likelihood estimation (MLE). VMPs provide enhanced extrapolation capability compared to ProMP, as they handle via-points (including start and target positions) adaptation to points that lie out of the demonstrated distributions. In this paper, we leverage VMPs to learn the demonstrated motion styles of each keypoint and to adapt the corresponding trajectories to via-points identified using the dense correspondence function  $f_c$  of DON. In contrast to control policies based on RL (e.g., [16]) or on visual servoing (e.g., [12]), VMP-based control policies endow K-VIL with flexible temporal scaling and reliable via-point adaptation.

## IV. KEYPOINT-BASED VISUAL IMITATION LEARNING

In this section, we present the proposed K-VIL approach. Given  $N$  demonstration videos  $\mathcal{V} = \{V_n\}_{n=1}^N$  of a task in  $D$ -dimensional task space, where  $D \in \{2, 3\}$ , K-VIL first preprocesses the RGB-D videos and generates the data required for learning the task. This includes densely sampled candidate points, their descriptors and trajectories, the spatial properties and roles of the objects, as well as all potential local frames (see Section IV-A). A sparse set of keypoints and their geometric constraints are then estimated via *principal constraint estimation* (PCE). As detailed in Section IV-B, our proposed PCE first extracts a set of keypoints and their geometric constraints by leveraging PME algorithms. These algorithms rely on observed distances, as well as on the demonstration variability when several demonstrations are provided. For the cases where the resulting set contains redundant selections of keypoints, our PCE then leverages Hierarchical Agglomerative Clustering (HAC) to resolve this redundancy and obtain a final set  $\mathcal{P}$  of sparse keypoints. As explained in Section IV-C,  $\mathcal{P}$  is then used to extract the task representation consisting of a set of keypoints defined by visual descriptors  $\mathcal{D} = \{d_i\}_{i=1}^L$ , their associated geometric constraints  $\mathcal{C} = \{C_i\}_{i=1}^L$  and the weights of the via-point movement primitives  $\Omega = \{w_i\}_{i=1}^L$ , which are then exploited to reproduce the keypoint motions. The extracted task representation is finally used by the keypoint-based admittance controller (KAC) presented in Section V to reproduce the demonstrated skill on the robot. The proposed K-VIL approach is shown in Fig. 2, and its different steps are detailed next. The main notations are listed in Table I.

### A. Preprocessing

As previously mentioned, K-VIL first preprocesses the RGB-D videos  $\mathcal{V}$  provided as demonstrations. This is achieved via the following five steps, also depicted in Fig. 2.

1) *Sampling of candidates*: First, we query the list of objects categories  $\mathcal{O} = \{O_i\}_{i=1}^I$  involved in the task by feeding the Mask R-CNN model with an image randomly sampled from  $\mathcal{V}$ . From the *visible region* of each object  $O_i$ ,  $P_i$  candidate points are then densely and uniformly sampled and form a set  $\mathcal{P}_i$ . Each candidate point is a potential keypoint or a potential origin of a local frame, and may later be selected as such by K-VIL. Note that we here assume that all relevant points are



Notation	Meaning	Notation	Meaning
$C$	geometric constraint	$\xi_1, \xi_2$	lower and upper thresholds of spatial variability
$d$	the intrinsic dimension of a principal manifold	$\lambda$	the regularization factor of PME
$D, \bar{D}$	the dimension of the task space, descriptor space	$\kappa_f$	the curvature of the principal manifold
$g_1, g_2$	force scaling parameters	$\mathbf{d}$	the descriptor vector of a keypoint / candidate
$H$	the total number of candidate points	$\mathbf{f}$	a force vector
$I$	the number of objects	$\mathbf{h}_c$	the Coriolis and gravitational force in task space
$L$	the number of constraints	$\mathbf{k}, \mathbf{k}$	the position and velocity vector of a keypoint
$N$	the number of demonstrations	$\mathbf{p}$	a position vector of a point
$P$	the number of points sampled on an object	$\mathbf{w}, \Sigma_w$	a weights vector of the VMP and its covariance
$Q$	the number of neighboring points	$\mu_w$	the mean of $\mathbf{w}$
$\mathcal{C}$	the set of constraints	$\nu$	the explained variance
$\mathcal{D}$	the set of descriptors	$\eta$	the spatial variability
$\mathcal{F}$	a local frame	$\mathbf{K}$	diagonal stiffness, damping and inertia matrices
$\hat{\theta}_{\mathcal{F}}$	the configuration of a canonical local frame	$\mathbf{S}$	the canonical shape
$\hat{\Theta}_{\mathcal{F}}$	the set of canonical local frame configurations	$\tau$	a trajectory
$\mathcal{M}$	a manifold	$\sigma(\cdot)$	the density force on the $d$ -dimensional manifold
$\mathcal{O}, \mathcal{O}$	an object category, the set of objects	$\nabla\sigma(\cdot)$	the density field
$\mathcal{P}$	the set of keypoints / candidate points	$f_c(\cdot, \cdot)$	the correspondence function
$\gamma, \mathcal{R}$	the role of the object, the set of object roles	$f(\cdot)$	the projection index of a principal manifold
$\mathcal{S}$	the set of canonical shapes of all objects	$\pi_d(\cdot)$	the reconstruction function of a principal manifold
$\mathcal{V}$	the set of demonstration videos	$\psi(\cdot)$	the squared exponential (SE) kernels in VMP
$\varphi, \Phi$	the spatial scale, the set of spatial scales	$\psi_i(\cdot)$	the squared exponential kernel

TABLE I: Summary of K-VIL's notations.

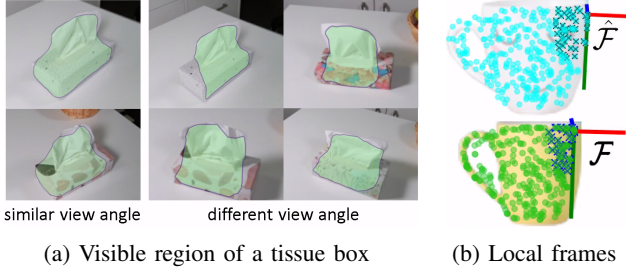


Fig. 3: (a) Illustration of the visible region of a tissue box. (b) Example of local frame matching on two cups. The canonical shape of the cup category (*top*) is defined on a white cup by the positions of its candidate points ( $\bullet$ ). The correspondence points ( $\bullet$ ) of the canonical shape are detected on a yellow cup (*bottom*). A canonical local frame  $\hat{\mathcal{F}}$  is parametrized by  $Q = 50$  neighboring candidates ( $\times$ ). These candidates are then used to find the same local frame  $\mathcal{F}$  on the yellow cup.

located in the region on the object surface that is always visible to the imitator (see Fig. 3a). We denote the set of all candidate points from all objects as  $\mathcal{P}_c = \bigcup_{i=1}^I \mathcal{P}_i$ . Their corresponding deep visual feature descriptors are derived from DON [5] as  $\mathcal{D}_c = \{\mathbf{d}_h\}_{h=1}^H$  with  $\mathbf{d}_h \in \mathbb{R}^3$  and  $H = |\mathcal{P}_c|$  the cardinality of  $\mathcal{P}_c$ .

2) *Trajectories of candidate points*: We extract the task space trajectory of all candidate points from the videos using the DON-based correspondence function  $f_c(\cdot, \cdot)$  (see Section III-A), which finds the correspondence pixel of the candidates and maps them to 3D coordinates in the camera local frame. The obtained trajectories are then smoothed and normalized in time with  $T$  timesteps. We obtain a set  $\mathcal{T} = \{\tau_h\}_{h=1}^H$  of trajectories of all candidates, where  $\tau_h \in \mathbb{R}^{N \times T \times D}$  denotes the trajectory of the  $h^{\text{th}}$  candidate point. These trajectories are used in the remaining preprocessing steps and in Section IV-B to extract keypoints and geometric constraints.

3) *Object properties*: We define the *canonical shape*  $\mathcal{S}_i \in \mathbb{R}^{P_i \times D}$  of each object category  $\mathcal{O}_i$  as the positions of

all candidates on the object at the first time step of the first demonstration. This notion of canonical shape is illustrated for a cup in Fig. 3b. Moreover, we define the *spatial scale*  $\varphi_i \in \mathbb{R}$  as the maximum distance between each pair of candidates on the canonical shape, which will be used in Section IV-B2 to determine object-independent thresholds.

4) *Object roles*: Geometric constraints do not suffice to entirely represent a task. For example, one of the constraints of a pouring task is the kettle-cup alignment, which could be achieved by moving the cup toward a static kettle. Instead, pouring requires a motion of the kettle. K-VIL addresses this issue by considering the role of the objects for the task at hand. Namely, we detect object motion saliency similarly to [29] to determine the role of the objects  $\mathcal{R} = \{\gamma_i\}_{i=1}^I$ , where  $\gamma_i \in \{\text{master}, \text{slave}\}$ . The master  $\mathcal{O}_m$  is the object with the lowest average variance of candidates' trajectory, while other objects are slaves  $\mathcal{O}_s$ . K-VIL accounts for the objects' roles by constructing local frames only on the master and extracting keypoints only on the slaves. Therefore, we split the set  $\mathcal{P}_c$  of all candidates by the object roles to  $\mathcal{P}_m$  and  $\mathcal{P}_s$ , denoting the set of candidates on master and slave objects respectively. Similarly,  $\mathcal{D}_c$  is splitted to  $\mathcal{D}_m$  and  $\mathcal{D}_s$ .

5) *Local frame detection*: As previously mentioned, K-VIL aims at representing the demonstrated task from an object-centric perspective. This is achieved by defining local frames on the master object. To do so, we initially define one canonical local frame  $\hat{\mathcal{F}}_j$  equal to identity for each candidate point  $j \in \mathcal{P}_m$  of the master canonical shape (see Fig. 3b-top). Each local frame  $\hat{\mathcal{F}}_j$  is assigned the  $Q$  closest candidates to  $j$ , whose positions in  $\hat{\mathcal{F}}_j$  are denoted as the reference values  $\mathbf{p}_q^*$ . In other words,  $\hat{\mathcal{F}}_j$  is parameterized by  $\hat{\theta}_{\hat{\mathcal{F}}_j} = \{\{\mathbf{d}_q\}_{q=1}^Q, \{\mathbf{p}_q^*\}_{q=1}^Q\}$ , where  $\mathbf{d}_q$  are the descriptors of the  $Q$  neighboring candidates. These neighboring candidates are then used to detect the same local frame on another instance of the same object category at a different time  $t$  (see Fig. 3b-bottom). Namely, the local frame  $\mathcal{F}_j$  is detected by minimizing the mean squared displacement of the observed coordinates  $\{\mathbf{p}_q(t)\}_{q=1}^Q$  of the neighboring

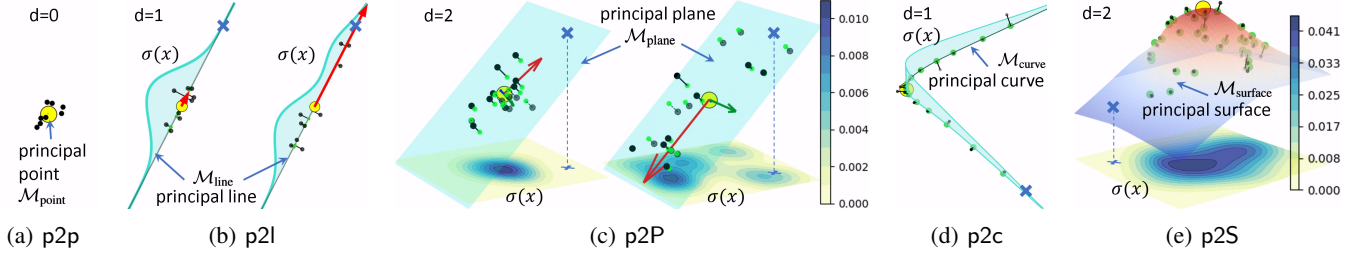


Fig. 4: Five types of geometric constraints. The constraints are obtained from candidate points (●) from  $N$  demonstrations. The density function  $\sigma(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^d$  is estimated from the projections (●) of the candidate positions on the  $d$ -dimensional principal manifold. We also depict the mean  $\mathbf{p}_m$  (●), the spatial variability ( $\rightarrow$ ,  $\rightarrow$ ) on the principal manifold along the principal components, the stress vector  $\mathbf{s}$  ( $\rightarrow$ ), and examples of extrapolation of the keypoints (×) on the manifolds.

candidates at time  $t$  with their reference values  $\{\mathbf{p}_q^*\}_{q=1}^Q$ , i.e.,

$$\mathcal{F}_j(t) = \arg \min_{\mathcal{F}} \sum_{q=1}^Q \|\mathbf{p}_q^* - \mathbf{p}_q(t)\|^2.$$

The set of local frames on the master object is denoted as  $\hat{\Theta}_{\mathcal{F}} = \{\hat{\vartheta}_{\mathcal{F}_j} : j \in \mathcal{P}_m\}$ .

In summary, the preprocessed data for all objects contain the sets  $\mathcal{P}_m$  and  $\mathcal{P}_s$  of the candidate points on master and slave objects, respectively, their corresponding descriptors  $\mathcal{D}_m, \mathcal{D}_s$  and trajectories  $\mathcal{T}$ , the set  $\mathcal{S} = \{\mathcal{S}_i\}_{i=1}^I$  of the object canonical shapes, the set  $\Phi = \{\varphi_i\}_{i=1}^I$  of the object spatial scales, the set  $\mathcal{R}$  of the object roles, and the set  $\hat{\Theta}_{\mathcal{F}}$  of all local frames.

### B. Principal Constraints Estimation (PCE)

Given the preprocessed data, our goal is to jointly extract a set  $\mathcal{P}$  of keypoints and a set  $\mathcal{C} = \{C_l\}_{l=1}^L$  of geometric constraints. As shown in Fig. 4, we consider five basic types of geometric constraints for keypoints in a 3D Cartesian space, namely point-to-point (p2p), point-to-line (p2l), point-to-plane (p2P), point-to-curve (p2c) and point-to-surface (p2S). The p2p, p2l, and p2P constraints are linear and can therefore be estimated by analyzing the variance of the keypoint positions in multiple demonstrations using PCA [38]. In contrast, p2c and p2S are nonlinear geometric constraints, which we estimate with the iterative PME (see Section III-B). Note that more complex constraints such as colinear, coplanar, parallel and perpendicular result from combinations of our five basic types of constraints. To ensure that the constraints are reliably estimated, the criterion of the proposed PCE is adapted to the number of demonstrations. Specifically, the single-demonstration case (i.e., one-shot IL) is considered a special case as it does not provide sufficient information to learn generalizable skills. Therefore, we propose heuristically-designed distance-based criteria (Section IV-B1). In contrast, when several demonstrations are available (i.e., few-shot IL), the keypoints and geometric constraints are learned based on the variability of the demonstrations (Sections IV-B2 and IV-B3). Moreover, nonlinear constraints are considered only if enough demonstrations are available.

1) *Distance criteria for a single demonstration:* A single demonstration ( $N = 1$ ) does not provide examples of variations in the demonstrated task, and thus prevents the learning of generalizable skills. Therefore, in this case, we learn a set of constraints that fully determines the pose of the objects. To do so, we assume that the objects are rigid and extract

3 keypoints for each slave object in order to fully constrain their position in a 3-D space. Note that K-VIL can also be applied to 2-D cases, where 2 keypoints are sufficient to determine the pose of an object. We map the trajectories of all candidates on the slave objects into each of the canonical local frames  $\mathcal{F}_j(t)$  on the master object at time step  $t$ , where  $j \in \mathcal{P}_m$ . Therefore, all demonstrations obtained from arbitrary viewpoints are aligned in a common viewpoint defined by the local frame  $\mathcal{F}_j(t)$  (see Section VI-B2). Then,  $\tilde{\tau}_j^k(t) \in \mathbb{R}^{N \times D}$  with  $k \in \mathcal{P}_s$  represents the positions of the  $k^{\text{th}}$  candidate point in all demonstrations viewed from the  $j^{\text{th}}$  common viewpoint at time step  $t$ , we use this variable to denote the *candidate positions* in the remaining of the paper. We observed that, for a variety of daily manipulation tasks, the closest point  $k_1$  on the slave object to the master object is often crucial to respect contact or avoid a collision, whereas the furthest point  $k_2$ , in combination with  $k_1$ , controls the pose of the object. Motivated by these heuristics, we propose the following procedure for each slave object. First, we choose the local frame  $\mathcal{F}^*(t)$  from the canonical local frames as the closest on average to all candidates on the slave objects. The two keypoints  $k_1, k_2$  on the slave object then correspond to the closest and farthest candidates from the selected local frame  $\mathcal{F}^*(t)$ . For a 3-D task space, we select an additional keypoint  $k_3$  as the farthest candidate from both  $k_1$  and  $k_2$ . To fully determine the pose of the slave object, the three keypoints are subject to linear p2p constraints. For each slave object, we finally obtain a set  $\mathcal{P}_d = \{k_l\}_{l=1}^D$  of keypoints and the corresponding geometric constraints  $\mathcal{C} = \{C_l\}_{l=1}^D$ , where  $C_l = \{\mathcal{M}_{\text{point}}(k_l), t, \theta_{\mathcal{F}^*(t)}, \text{p2p}\}$  defines a p2p constraint on a 0-dimensional principal manifold  $\mathcal{M}_{\text{point}}$  on point  $k_l$  at time  $t$  in the local frame given by  $\vartheta_{\mathcal{F}^*}$ .

2) *Variance criteria for linear constraints:* When several demonstrations ( $N > 1$ ) are available, we leverage their variability to estimate linear constraints beyond p2p. To do so, we obtain the candidate positions  $\tilde{\tau}_j^k(t)$  in the canonical local frames  $\mathcal{F}_j(t)$  at time  $t$  as described in Section IV-B1 and compute the explained variance  $\nu_j^k(t) = \mathbb{V}_{\text{PCA}}[\tilde{\tau}_j^k(t)] \in \mathbb{R}^D$  of each candidate using PCA. The *spatial variability*  $\eta_j^k(t)$  is then defined as  $\eta_j^k(t) = (\nu_j^k(t))^{1/2} / \tilde{\varphi}_i$ , with  $\tilde{\varphi}_i$  the spatial scale of the slave object  $O_i$  to which the  $k^{\text{th}}$  candidate belongs. In contrast to the explained variance, spatial variability removes dependencies on the object size. This allows us to empirically define two object-agnostic lower and upper thresholds  $\xi_1, \xi_2$  to identify appropriate linear geometric constraints based on the

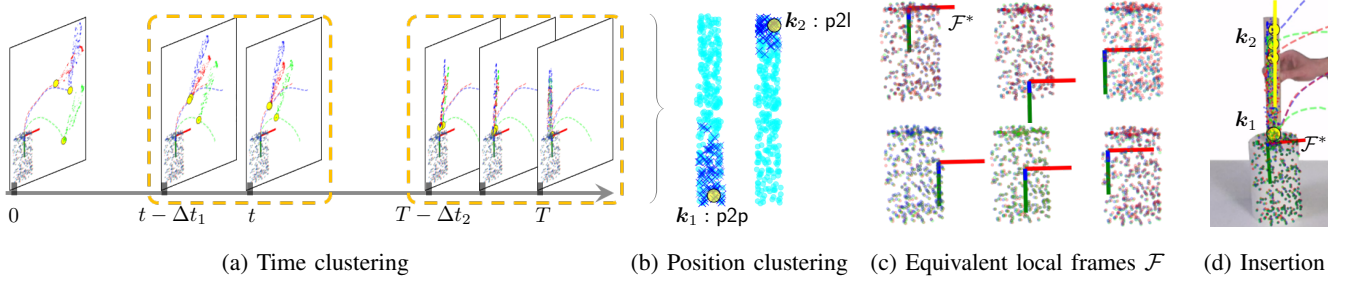


Fig. 5: Hierarchical agglomerative clustering for inserting a stick into a paper roll (see also Section VI for the task description). (a) The selected candidates are first clustered in time (---) to identify the adjacent timesteps. (b) The candidates ( $\times$ ) in each time cluster (e.g., here in the last time cluster  $[T - \Delta t_2, T]$ ) are then clustered based on their positions in the canonical shape ( $\bullet$ ) of the stick for each constraint (here p2p and p2l). For each position cluster, the keypoint ( $\bullet$ ) with the lowest variability is finally selected. (c) Since the paper roll has no shape variation, i.e., all canonical local frames are equivalent, the closest frame  $\mathcal{F}^*$  to the selected keypoints is selected. (d) Final task representation (see also Fig. 13).

computed spatial variability. Namely, the constraints of each candidate  $k$  are determined by the following three conditions: (i)  $\eta_{j,1}^k(t) < \xi_1$  implies a low spatial variability in the first components, as well as across all other dimensions since all the components of the spatial variability are ranked in decreasing order, i.e.,  $\eta_{j,e}^k > \eta_{j,e+1}^k$ , where  $e = \{1, 2\}$ . This also means that the position of the candidate  $k$  remains close to a fixed point  $\mathcal{M}_{\text{point}}$  across all demonstrations. Therefore,  $k$  is subject to a p2p constraint; (ii)  $\eta_{j,2}^k(t) < \xi_1$  and  $\eta_{j,1}^k(t) > \xi_2$  imply that  $k$  is constrained on a line  $\mathcal{M}_{\text{line}}$  along the first component, i.e.,  $k$  is subject to p2l constraint; (iii) Similarly,  $\eta_{j,3}^k(t) < \xi_1$  and  $\eta_{j,2}^k(t) > \xi_2$  indicate that  $k$  is constrained on a plane  $\mathcal{M}_{\text{plane}}$  going through the first two components, i.e.,  $k$  is subject to a p2P constraint. Any spatial variability  $\eta_j^k(t)$  satisfying the above conditions indicates the joint selection of the  $k^{\text{th}}$  candidate, the  $j^{\text{th}}$  local frame and  $t^{\text{th}}$  time step. All candidates selected as such form a set  $\mathcal{P}_l$  of keypoints subject to linear geometric constraints. Note that, due to the fact that two distinct points define a line and three non-collinear points define a plane, we learn p2l constraints when  $N > 2$  and of p2P constraints when  $N > 3$ .

3) *Variance criteria for nonlinear constraints*: The linear constraints may not suffice to represent a given task accurately despite them being easily estimated from a few demonstrations. For instance, the pouring task of Fig. 1c requires a point-to-curve constraint. Therefore, we additionally estimate nonlinear (p2c and p2S) constraints with PME. In the following, a set  $\mathcal{P}_s$  of all candidate points on slave objects that do not satisfy any linear constraints are considered as potential candidates for nonlinear constraints. In our case, we replace the random  $D$ -dimensional vector  $\mathbf{x}$  in (1) with the candidate point  $\tilde{\tau}_j^k(t)$  on the demonstrated trajectory  $\tilde{\tau}_j^k$  at time step  $t$ , so that the PME loss in Section III-B becomes

$$\mathcal{L}(f, \pi_d) = \mathbb{E} \left\| \tilde{\tau}_j^k(t) - f(\pi_d(\tilde{\tau}_j^k(t))) \right\|^2 + \lambda \|\kappa_f\|^2, \quad k \in \mathcal{P}_s.$$

After obtaining  $\pi_d$  from PME, we compute the projections of candidates onto the manifold, i.e.,  $\hat{\tau}_j^k(t) = \pi_d(\tilde{\tau}_j^k(t))$ , where  $\hat{\tau}_j^k(t) \in \mathbb{R}^{N \times d}$ . Then, analogously to Section IV-B2, we define the explained variance  $\nu_{j,\parallel}^k(t)$  in the tangential direction of the principal manifold as the variance of the projections, i.e.,  $\nu_{j,\parallel}^k(t) = \mathbb{V}[\|\hat{\tau}_j^k(t)\|] \in \mathbb{R}$ . The explained variance  $\nu_{j,\perp}^k(t)$  in the

orthogonal direction corresponds to the variance of the length of the stress vectors  $\mathbf{s} = \hat{\tau}_j^k(t) - f(\hat{\tau}_j^k(t))$ , i.e.,  $\nu_{j,\perp}^k(t) = \mathbb{V}[\|\mathbf{s}\|]$ . Similar to the linear case, the spatial variability is defined as  $\eta_{j,z}^k(t) = \sqrt{\nu_{j,z}^k(t)}/\bar{\varphi}_i$ ,  $z \in \{\perp, \parallel\}$ . The set  $\mathcal{P}_{nl}$  of keypoints subject to nonlinear geometric constraints is then selected as

$$\mathcal{P}_{nl} = \{k \mid \nu_{j,\perp}^k(t) < \xi_1, \nu_{j,\parallel}^k(t) > \xi_2, k \in \mathcal{P}_s\}.$$

The type of the geometric constraints is determined by the intrinsic dimension  $d$  of the learned principal manifold, i.e.,  $d = 1$  and  $d = 2$  indicate a p2c and a p2S constraint, respectively. Notice that, in order to guarantee their reliable estimation, nonlinear constraints are considered only when enough demonstrations ( $N > 10$ ) are available.

4) *Hierarchical Agglomerative Clustering (HAC)*: As explained in Section IV-B2, each selected candidate point  $k$  in the resulting sets of linear and nonlinear constraints  $\mathcal{P}_l$  and  $\mathcal{P}_{nl}$  corresponds to jointly selected time step  $t$  and local frame  $j$ . Redundancy may occur due to adjacent timesteps, neighboring keypoints, or equivalent local frames. To resolve this redundancy, we first cluster the keypoints in time using Hierarchical Agglomerative Clustering (HAC) to identify adjacent timesteps. Fig. 5 shows an example of HAC for an insertion task. We then use HAC again to cluster the keypoints within each time cluster based on their positions in the canonical shape of the slave object, thus identifying neighboring keypoints. The redundancy is finally resolved by keeping only the keypoint with the lowest variability to represent each position cluster. This keypoint is selected for its robustness against sensor and correspondence detection noise. If a selected keypoint at a selected time step is subject to multiple constraints represented in different local frames, we select the closest local frame to the keypoint on average. In summary, the proposed PCE retrieves a sparse set of  $L$  keypoints as the union  $\mathcal{P} = \mathcal{P}_d \cup \mathcal{P}_l \cup \mathcal{P}_{nl}$  and their associated (non)linear constraints  $\mathcal{C} = \{\mathcal{C}_l\}_{l=1}^L$ , which are exploited to represent the task as explained next.

### C. Extraction of $K$ -VIL's complete task representation

While the keypoints and associated constraints estimated in Section IV-B allow us to understand the demonstrated task, a



control policy is additionally required for reproducing the task. Here, we propose to model the observed keypoints trajectories as VMPs [13]. For our purposes, we train the VMPs from an object-centric perspective and according to the constraints estimated via PCE. Specifically, for each keypoint subject to a p2p constraint, a VMP is trained on its observed trajectory  $\tilde{\tau}_l$  retrieved in the corresponding local frame  $j$  from time step 1 to the extracted time step  $t$ , i.e.,  $\tilde{\tau}_l = (\tau_j^k(1) \dots \tau_j^k(t))^T$ , where  $k$  and  $l$  indicates that the  $k^{\text{th}}$  candidate point in the dense set  $\mathcal{P}_s$  is selected as the  $l^{\text{th}}$  keypoint in the sparse set  $\mathcal{P}$ . Note that, for the case of intrinsic dimension  $d > 0$  (i.e., p2l, p2P, p2c, p2S constraints), the constraint is fulfilled if and only if the corresponding keypoint is placed at the time step  $t$  on the principal manifold that defines the constraint. Although the location of the keypoint on the manifold does not affect the fulfillment of the constraints, it may influence the similarity between the demonstrated object poses and those obtained in the reproduction. Therefore, we propose to decompose the control of such keypoint by considering the orthogonal and the tangential direction with respect to the corresponding principal manifold independently (see Fig. 6). The keypoint motion along the orthogonal direction represents the demonstrated style of approaching the principal manifold and guarantees the fulfillment of the constraints, while the motion along the tangential direction realizes the extrapolation of the keypoint target position and controls the similarity of the object pose between the demonstrations and the reproduction. Due to potential large shape variations in the objects used when reproducing the task, the final keypoints' target positions on the principal manifold may not align with the demonstrated targets. Therefore, we only train the VMP on the keypoint trajectories projected onto the orthogonal direction, i.e.,  $\tilde{\tau}_{l,\perp}$ . An example of the projected and reproduced trajectories obtained using the learned VMP in the case of a p2c constraint is shown in Fig. 7. At each time step during reproduction, we uncover the 3D target position of a keypoint by adding an offset generated by the VMP in the orthogonal direction to the orthogonal projection of the current keypoint onto the principal manifold (see  $k_2^*$  in Fig. 8 and Fig. 9). Therefore, by setting the VMP goal to 0, the keypoints fulfill the corresponding geometric constraints at the end of their trajectory. The motion of the keypoints along the orthogonal and tangential directions is controlled via the keypoint-based admittance controller presented in the next section. In summary, K-VIL's final task representation is composed of a set of keypoints represented by their descriptors  $\mathcal{D} = \{d_l\}_{l=1}^L$ , their associated geometric constraints  $\mathcal{C} = \{C_l\}_{l=1}^L$ , and their associated movement primitives encoded via the set of weights  $\Omega = \{w_l\}_{l=1}^L$ .

## V. KEYPOINT-BASED ADMITTANCE CONTROLLER

After learning the representation of a given task from demonstrations, we aim at reproducing this task with a robot. This means that the robot should be able to interact with the objects such that their keypoints follow the learned constrained trajectories. This requires filling the gap between K-VIL's task representation and real-time robot controllers. To this end, we propose a Keypoint-based Admittance Controller (KAC),

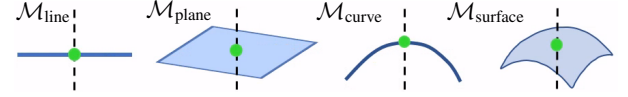


Fig. 6: Orthogonal direction (---) to the principal manifolds.

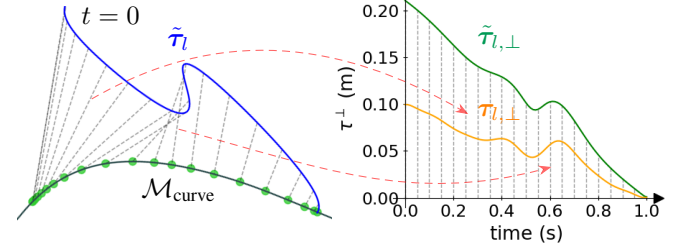


Fig. 7: Projected and reproduced trajectories using a VMP for a p2c constraint. The demonstrated trajectory  $\tilde{\tau}_l$  (—) is projected at each time step in the orthogonal direction (---) of the principal manifold  $\mathcal{M}_{\text{curve}}$ . The projected trajectory  $\tilde{\tau}_{l,\perp}$  (—) is used to train movement primitives, which is then used to reproduce trajectories, e.g.,  $\tau_{l,\perp}$  (—) with a new start position at 0.1 m and goal position at 0 m. The arrows (→) mark the corresponding projected trajectory between the 3D and 2D plots at two timesteps.

which (i) handles variable numbers of extracted keypoints for different tasks; (ii) enables the extrapolation of keypoint target positions on their learned principal manifolds; (iii) resolves potential interference between different types of geometric constraints. Note that (ii) and (iii) are required to handle large object shape variations in the task reproduction.

Specifically, a KAC associates each keypoint in  $\mathcal{P}$  with a virtual spring-damper system, whose attractor is computed via the corresponding VMPs (see Section V-A). As detailed in Section V-D, the sum of the attraction forces of the spring-damper systems of all keypoints is then used as the task-space force command for the robot. This allows the KAC to handle a varying number of keypoints for different tasks. Regarding (ii), the extrapolation of keypoint target positions subject to a p2p constraint is not allowed. For non-p2p constraints, this is achieved by decomposing the control in orthogonal and tangential directions of the learned principal manifolds (see Section IV-C). As a result, the keypoints approach the principal manifolds using the motion profiles learned from the projected trajectories in orthogonal directions. The control force generated by the virtual spring-damper system of each keypoint remains orthogonal to the principal manifold at each time step, and the keypoints reach the corresponding geometric constraints when the execution of the VMP finishes. While this leads to the extrapolation of keypoint target positions on the principal manifold, it does not account for the distance between the demonstrated and extrapolated targets. Therefore, we propose to balance extrapolation and regulation by estimating the density function of the demonstrated targets on the principal manifolds, as described in Section V-B. This density is then used to compute an additional force, i.e., the density force, that drives each keypoint toward the demonstrated targets. Finally, we address the interference issue (iii) by assigning different priorities to different types of geometric constraints in Section V-C. The different steps of KAC are

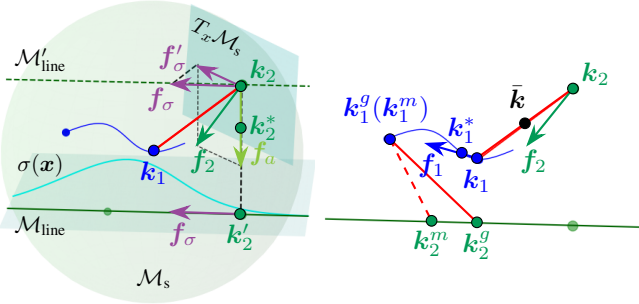


Fig. 8: Illustration of the attraction and density forces when the keypoints  $k_1$  and  $k_2$  are subject to p2p and p2l constraints, respectively. **Left:** The approach force  $f_a$  of  $k_2$  is computed by the virtual spring-damper system between the attractor  $k_2^*$  and  $k_2$ . The density force  $f_\sigma$  is then projected onto the tangent space of the sphere at  $k_2$ . The control force of  $k_2$  is the combination of the attraction and the projected density force. **Right:**  $k_1$  is controlled by the attraction force  $f_1$  following the attractor  $k_1^*$  and the VMP (—) to reach the target  $k_1^g$ , which coincides with the demonstrated target  $k_1^m$ . Note that the target  $k_2^g$  of  $k_2$  does not coincide with  $k_2^m$  due to object shape variation, i.e., the distance (—) between  $k_1$  and  $k_2$  during reproduction is longer than for the demonstration (---).

detailed next.

#### A. Attraction force

Given K-VIL's task representation and a new image frame  $A$  for the task reproduction, we can identify the keypoints representing the task. Namely, their positions  $k_l \in \mathbb{R}^D$ , with  $l \in [1, L]$ , represented in the root frame  $\mathcal{F}_r$  of the robot are obtained using the visual descriptors  $d_l$  and the DON-based correspondence function  $f_c(A, d_l)$  (see Section III-A). The attractor  $k_l^*$  of the virtual spring-damper system at each time step is computed for each keypoint by the corresponding VMPs projected onto  $\mathcal{F}_r$ . The attraction force generated by the virtual spring-damper system is then computed as

$$f_{a,l} = \bar{K}_p(k_l^* - k_l) + \bar{K}_d(\dot{k}_l^* - \dot{k}_l),$$

where  $\bar{K}_p, \bar{K}_d$  are diagonal stiffness and damping matrices, respectively, and  $\dot{k}_k$  and  $\dot{k}_l^*$  are the velocity of  $k_l$  and  $k_l^*$ , respectively. As explained in Section IV-C, in the case of non-p2p constraints, the VMPs are trained on trajectories projected in directions that are orthogonal to the principal manifold. Therefore, the learned VMPs and the attraction forces enable the reproduction of the demonstrated motion patterns in the orthogonal direction. This implies that the final positions of the keypoints can be extrapolated anywhere on the principal manifolds, e.g., to satisfy object shape variations and other geometric constraints. However, without considering the demonstrated targets on the principal manifold, we may lose important information about successful task execution or a specific style of execution. We obtain such information using kernel density estimation and provide additional density forces driving the keypoints toward the demonstrated targets.

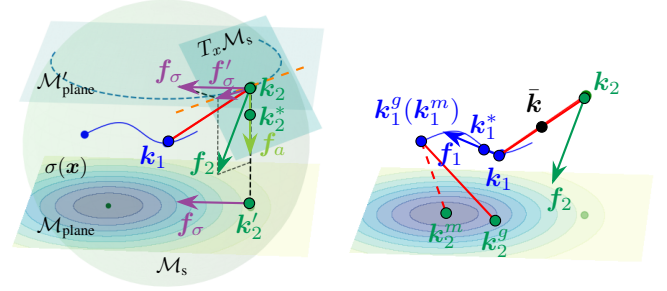


Fig. 9: Illustration of the attraction and density forces when  $k_2$  is subject to a p2P constraint. In contrast to Fig. 8, the density force  $f_\sigma$  is projected onto the intersection line (---) of the shifted principal manifold  $\mathcal{M}'_{\text{plane}}$  and the tangent space. Legend as in Fig. 8.

#### B. Density force

Given a non-p2p constraint, we project the demonstrated keypoint positions  $\tilde{\tau}_l(t)$  at the extracted time step  $t$  onto the corresponding  $d$ -dimensional principal manifold using the learned projection index  $\pi_d$ , i.e.,  $\hat{k}_l^m = \pi_d(\tilde{\tau}_l(t)) \in \mathbb{R}^{N \times d}$ . Since at time step  $t$ , the  $l^{\text{th}}$  keypoint is supposed to fulfill the geometric constraints, we interpret  $\hat{k}_l^m$  as its demonstrated target positions on the manifold. We then estimate the density function  $\sigma(x)$  of the keypoint target positions from  $\hat{k}_l^m$  using kernel density estimation [39] with SE kernels. Examples of estimated density functions for p2l, p2P, p2c, and p2S constraints are depicted in Figs. 4b to 4e. This density function indicates the probability of a keypoint target position on the corresponding principal manifold given the demonstrated target positions. In other words, the density function indicates the confidence level of K-VIL when extrapolating the keypoint target positions to new locations on the principal manifold, which may occur due to object shape variations in the reproduction. Examples of extrapolated keypoint target positions ( $\times$ ) during reproduction are depicted in Figs. 4b to 4e. Notice that the target position in Fig. 4b-left has a lower probability (i.e., a lower extrapolation confidence) than the one in Fig. 4b-right due to its increased distance with the demonstrated target positions. This illustrates that the control of such keypoints must not only fulfill the geometric constraints, but also be as close as possible to the demonstrated targets on the constraints. Therefore, in addition to the attraction force  $f_{a,l}$  that guarantees the fulfillment of the geometric constraint, we define a *density force*  $f_{\sigma,l}$  to drive the keypoints into regions with higher probability. To do so, we first define the driving force  $f_{\sigma,1}$  computed from the density field  $\nabla\sigma(x)$  as

$$f_{\sigma,1} = g_1 \cdot f(\nabla\sigma(x)). \quad (2)$$

For the regions where  $f_{\sigma,1}$  is too small to drive the keypoints, we then define a minimal driving force  $f_{\sigma,2}$  as

$$f_{\sigma,2} = g_2 \cdot f\left(\frac{k_l^m - k_l^i}{\|k_l^m - k_l^i\|_2}\right), \quad (3)$$

where  $k_l^i = \pi_d(k_l) \in \mathbb{R}^d$  is the projection of the keypoint onto the principal manifold,  $f(\cdot)$  is the reconstruction function (see Section IV-B3),  $g_1$  and  $g_2$  are the force scaling parameters. Note that  $f_{\sigma,2}$  points directly to the mean of the demonstrated



targets  $\hat{k}_l^m$  on the principal manifold, i.e.,  $k^m = \text{avg}(\hat{k}_l^m)$ . Then the density force is the maximum of  $f_{\sigma,1}$  and  $f_{\sigma,2}$ , i.e.,

$$f_{\sigma,l} = \arg \max_f \|f\|_2, \quad f \in \{f_{\sigma,1}, f_{\sigma,2}\} \quad (4)$$

Examples of such density forces  $f_\sigma$  in the case of p2l and p2P constraints are shown in Figs. 8 and 9.

### C. Priority

In the case of large object shape variations, controlling a p2l constraint with the same priority as a p2p constraint may lead to a violation of the latter. To reduce such interference, we propose to set a higher priority to p2p constraints compared to the other constraint types. For the sake of clarity, we use Figs. 8 and 9 to explain this concept, where Fig. 8 shows the case of two constraints, p2p for  $k_1$  and p2l for  $k_2$ , while  $k_2$  is subject to a p2P constraint in Fig. 9. In both cases, we construct a sphere centered at  $k_1$  with radius  $\|k_2 - k_1\|$  and define the tangent space of the sphere at  $k_2$  as the plane formed by all the lines tangent to the sphere at  $k_2$ . For clarity, Figs. 8 and 9 also depicts the corresponding principal manifolds  $\mathcal{M}_{\text{line}}, \mathcal{M}_{\text{plane}}$  shifted in parallel to go through  $k_2$  as  $\mathcal{M}'_{\text{line}}, \mathcal{M}'_{\text{plane}}$ . Assuming solid connections (—) between  $k_1$  and  $k_2$ , large density forces  $f_\sigma$  generated for  $k_2$  will also drag  $k_1$  along the same direction. This may lead to the violation of the p2p constraint of  $k_1$ , and cause collision if  $k_1$  is close to the master object. To reduce such interference when the principal manifold is a line  $\mathcal{M}_{\text{line}}$ , we project  $f_\sigma$  onto the tangent space  $T_{x_s}\mathcal{M}_s$  of the sphere  $\mathcal{M}_s$ , so that the motion of  $k_1$  remains unaffected by the projected density force  $f'_\sigma$  (see Fig. 8). Similarly, when a principal manifold is a plane (Fig. 9), we project  $f_\sigma$  onto the intersection between the tangent space and the shifted principal plane  $\mathcal{M}'_{\text{plane}}$ . This also holds for the nonlinear constraints (p2c and p2S), for which a linear approximation is considered at each time step.

In summary, on the one hand, the density force allows the reproduced task to be similar to the demonstrations on the principal manifold. On the other hand, the priority mechanism reduces the interference of the density force with p2p constraints, while maintaining the extrapolation capability of K-VIL. Overall, the decomposition of the control force into attraction force (Section V-A) and projected density force (Sections V-B and V-C) is key to balancing the similarity of the reproduced task to the demonstration and the extrapolation capability. In practice, the stiffness and damping gains of the virtual spring-damper systems are empirically tuned for good tracking accuracy and control stability. Notice that one-shot IL is considered as a special case (see Section IV-B1). This is due to the fact that the learned task representation is composed of 3 keypoints subject to p2p constraints. In this case, no density force is needed, and the constraint priorities are defined as  $\text{Pri}_1 > \text{Pri}_2 > \text{Pri}_3$ , since  $k_1$  usually represents the contact point of two objects. Therefore, to ensure a higher control precision of  $k_1$ , we set the stiffness gains of the three keypoints to  $\bar{K}_{p,1} = 5\bar{K}_{p,2} = 10\bar{K}_{p,3}$  and the respective damping gains to  $\bar{K}_{d,l} = 2\bar{K}_{p,l}^{1/2}$ , where  $l \in [1, 3]$ , which ensure a critically damped behavior for control stability.



Fig. 10: Objects used in our paper include (a) tissue boxes, (b) teacups, (c) a rack and a hat, (d) kettles, (e) a paper roll, (f) sticks, (g) dustpans and brushes. Note that the rack can be assembled with stick #6-10 to have multiple shape variations.

### D. Admittance control

The goal of the KAC is to compute the control command of the robot arm from the attraction forces  $f_a$  and the projected density forces  $f'_\sigma$  of all keypoints. To do so, we first compute the control force of each keypoint as  $f_l = f_{a,l} + f'_{\sigma,l}$  and define a virtual tool-center-point (TCP)  $\bar{k} = \sum_{l=1}^L k_l / L$  as the mean of all keypoint positions (see Fig. 8-right, Fig. 9-right). The virtual TCP is driven by a virtual force and torque

$$f_f = \sum_{l=1}^L f_l \quad \text{and} \quad f_\tau = \sum_{l=1}^L (k_l - \bar{k}) \times f_l$$

with  $\times$  denoting the vector cross product. The total control force  $f_v = [f_f^\top, f_\tau^\top]^\top$  is applied to the robot end-effector (i.e., the humanoid hand) to calculate the virtual acceleration as

$$\ddot{x}_v = \tilde{K}_p(x_0 - x_v) - \tilde{K}_d\dot{x}_v - \tilde{K}_m f_v,$$

where  $x_0, x_v$  are the initial and virtual poses of the robot end-effector,  $\dot{x}_v$  is its virtual velocity, and  $\tilde{K}_m, \tilde{K}_d$ , and  $\tilde{K}_p$  are the inertia, damping and stiffness factors, respectively. The robot is controlled using a task space inverse dynamics controller, whose task space control force  $f_m$  is calculated as

$$f_m = K_p(x_v - x) + K_d(\dot{x}_v - \dot{x}) + h_c,$$

where  $x, \dot{x}$  are the current end-effector pose and velocity,  $K_d$ , and  $K_p$  are the damping and stiffness factors of the impedance controller, respectively, and  $h_c$  represents the Coriolis and gravitational force in the task space.

## VI. EVALUATION

We evaluate our approach in five daily tasks involving different types of geometric constraints and various categorical

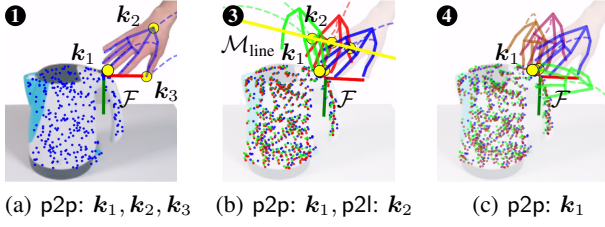


Fig. 11: Press a button on the kettle to open the lid with variations in hand orientation. The candidate points (colored points) and the skeleton of hands (colored line segments) are overlain on the objects at time step  $T$ . Different colors denote different trials. The keypoints  $k_l$  ( $\circ$ ) are extracted from different number  $N$  of demonstrations in each subfigure. The demonstrated trajectories ( $--$ ), the local frames  $\mathcal{F}$ , and the estimated principal manifolds ( $---$ ) are also depicted. Notice that the demonstrations were provided from different viewpoints, although they are here represented aligned to the local frame  $\mathcal{F}$  for a better illustration of the extracted keypoints and constraints (for viewpoint mismatch, see Fig. 16.)

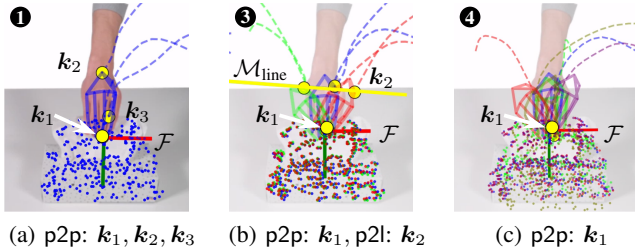


Fig. 12: Fetch tissue with variations in hand orientation. Legend as in Fig. 11.

objects (see Fig. 10). Namely, the considered tasks are press button (PB, Fig. 11), fetch tissue (FT, Fig. 12), insert sticks into a paper roll (IS, Fig. 13), pour water (PW, Fig. 14), hang hat on a rack (HH, Fig. 15), clean table with a dustpan and a brush (CT, Fig. 16). As a prerequisite for our experiments, we train DON in a self-supervised and task-agnostic manner. The training dataset was collected with a handheld Azure Kinect camera moving around objects such as tissue boxes, teacups, a rack, a hat, kettles, a paper roll, sticks, dustpans and brushes (see Fig. 10). The collected data was then post-processed via a 3D reconstruction process using Open3D [40]. We used MediaPipe to detect 21 keypoints on the human hands and treat the hands as a special type of object. It is important to emphasize that K-VIL is not limited to DON and MediaPipe, but instead can be used with any correspondence detection model, e.g., NDFs. In all experiments, we sample  $P_i = 300$  (for hands  $P_i = 21$ ) candidate points on each slave object and use  $Q = 50$  (for hands  $Q = 10$ ) neighboring candidates on the master object as references for local frame detection. We use 20 kernels for the VMPs. The empirical thresholds  $\xi_1, \xi_2$ , and the controller gains are fine-tuned for each task. A full list of control parameters is included in the example code.

We first evaluate the ability of K-VIL to extract generalizable task representations given different number of demonstrations (see Section VI-B1 and Section VI-B3 for a one-shot and a few-shot visual imitation learning setup, respectively, as

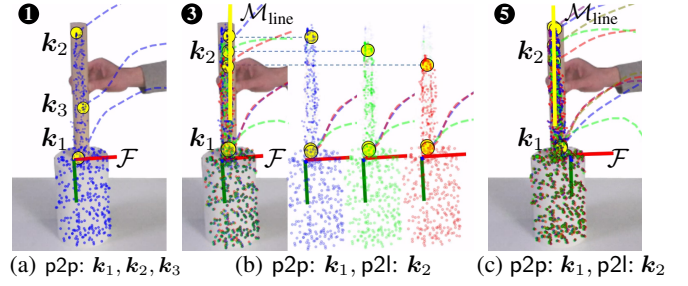


Fig. 13: Approach the insertion position to insert sticks with 3 length variations into a paper roll. Legend as Fig. 11.

well as Figs. 11 to 16). We demonstrate how variations in objects' pose and shape contribute to the efficient extraction of generalizable task representations and evaluate the ability of K-VIL to reproduce the corresponding tasks learned from a different number of demonstrations. We discuss the problems that arise when only scarce demonstrations are provided in Section VI-B4. We then show how they are resolved by providing more demonstrations and summarize the number of demonstrations required to learn a generalizable representation for each task. The proposed KAC is finally evaluated in terms of the control accuracy, precision, and success rate in Section VI-C. For more visualizations of the evaluation results on one/few-shot imitation learning, reproduction of skills learned from a different number of demonstrations, and other types of geometric constraints, we refer the interested reader to the accompanying videos and to the paper website.

#### A. Evaluation Protocols

For each task, we record a few demonstration videos (RGB-D) of a human performing the task using an Azure Kinect mounted on the head of the humanoid robot ARMAR-6 [41]. For tasks involving categorical objects, we distinguish between the object instances used for training of the vision models (the DON and Mask R-CNN models), for the demonstrations, and for the reproductions. If we only have one instance of a specific object category, this instance is used for the training, demonstrations, and reproductions. We define a set of *extraction tasks*  $\mathcal{T}_E = \{\text{PB, FT, PW, HH, IS, CT}\}$  for which we evaluate K-VIL's ability to extract generalizable task representations given  $N \in \{1, 3, 4, 5, 11\}$  demonstrations, respectively. For clarity, we only evaluate the representations of the last time cluster, i.e., the goal configuration of each task when  $t = T$ . We then define a set of *reproduction tasks*  $\mathcal{T}_R = \{\text{Task } N : \text{Task} \in \mathcal{T}_E\}$ , each of which is the reproduction of the Task by ARMAR-6 with the task representation extracted from  $N$  demonstrations. In order to evaluate the reproduction and adaptation of the learned task representations, e.g., the geometric constraints, in new cluttered scenes, the scene is perturbed arbitrarily before each trial of execution. Specifically, the involved objects and the robot hands are placed in arbitrarily different locations within the workspace and the view of the camera. The first image frame captured by the robot is used to parameterize the task with K-VIL's representation. This includes optimizing the local frames, identifying the keypoints, configuring the geometric



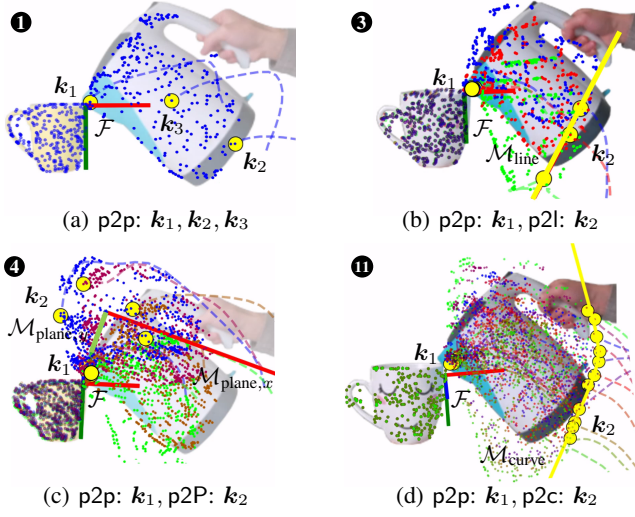


Fig. 14: Pouring task with variations in the shape of cups and the orientation of the kettle. The principal plane in (b) is represented by orthogonal vectors  $\mathcal{M}_{\text{plane},x}$  (—) and  $\mathcal{M}_{\text{plane},y}$  (—). Other legends as in Fig. 11.

constraints, and generating keypoint motion trajectories using the learned VMPs. We consider a task learned from  $N$  demonstrations and from a third-person view to be *generalizable* when it can be successfully reproduced by the robot with categorical objects in new cluttered scenes. Next, we describe the specifications of each considered task in terms of the collection of demonstrations and successful reproductions by the robot. Table VIII provides the list of considered tasks.

**Press Button (PB):** A human demonstrates how to open the lid of a kettle by pressing the corresponding button with the tip of the middle finger of either the left or the right hand. The kettle #5 of Fig. 10d is considered in this task. Both hands look similar and the demonstrator approaches the button with different hand poses (e.g., see Fig. 11). To reproduce the demonstrated human motion by the robot, we design fixed maps between keypoints of human hands and keypoints on the robot hand. The reproduction of the PB task is considered successful if the robot reaches the button with its fingertip within 5 mm to the target (the button) and if the lid is opened by closing the finger with a small angle.

**Fetch Tissue (FT):** A human demonstrates how to fetch tissue from two tissue boxes (#3 and #4 in Fig. 10a) with different hand poses (see Fig. 12). The tissue boxes #1-3 are used to train the vision models, whereas #4-10 are used for reproduction. The reproduction is considered successful if the robot can successfully grasp the tissue and pull it out of the boxes with a predefined pulling action. Although the shape variations between the tissue boxes #3-4 in the demonstrations are not obvious, box #10 introduces large shape variations for reproduction.

**Insert Stick (IS):** The sticks #2-4 in Fig. 10f are used to train the vision models and to demonstrate the insertion task by a human (see Fig. 13). Note that we do not insert the stick into the paper roll, as otherwise the keypoints will be occluded (we defer tasks with occlusion to future work). No pose variations are considered in this task. However, shape

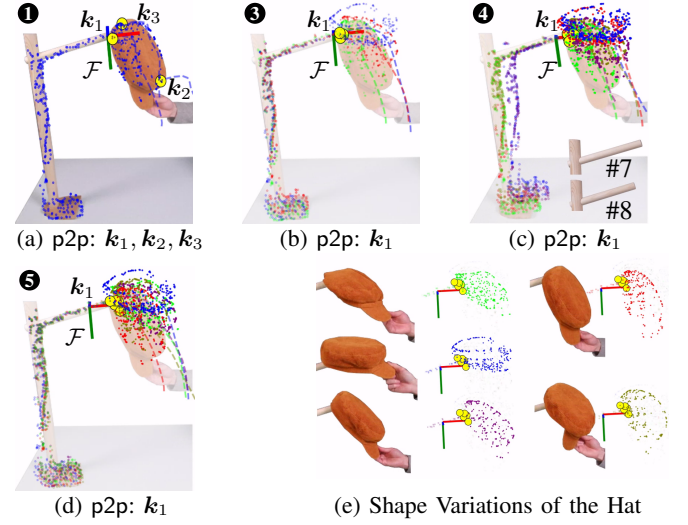


Fig. 15: Hang a hat on a rack. Note that there are only slightly deformations of the hat in (a)-(c) and relatively more obvious deformations in (d) (see (e)), while pose variations in the hats are considered in all cases. The racks in (c) have two shape variations. Legend as Fig. 11.

variations are introduced via sticks of different lengths and thicknesses. Moreover, we place the sticks with an initial tilting of  $\sim -30^\circ$  to  $80^\circ$  in the reproduction, thus extrapolating the demonstration range ( $\sim 0^\circ$  to  $45^\circ$ ). A successful reproduction is obtained by placing the lower tip of the stick right above the hole in the center of the paper roll without collision with the paper roll during execution (see Fig. 10e).

**Pour Water (PW):** The vision models for this task are trained with the teacups #5-8 in Fig. 10b and the kettles #4-5 in Fig. 10d. A human demonstrates the pouring task several times with the kettle #5 and teacups #1 and #3. The demonstrations incorporate teacup shape variations and kettle pose variations (see Fig. 14). The teacups #1-4 and all the kettles are used in reproduction. The reproduction of the PW task is successful if the spout of the kettle aligns above the rim of the teacup and the kettle is tilted appropriately.

**Hang Hat (HH):** The rack can be assembled with different lengths of sticks (#6-10 in Fig. 10f). The racks assembled with sticks #7-8 are used for training the vision models and for the demonstrations and #6-10 are used for the reproduction. In particular, the stick #7 is used in Figs. 15a, 15b and 15d and the sticks #7-8 are used in Fig. 15c. Successful reproductions are observed if the rim of the hat is placed on top of the tip of the stick regardless of the stick length and the initial pose of the hat.

**Clean Table (CT):** The dustpans #2-3 and the brushes #3-4 in Fig. 10g are used for training the vision models and for the demonstrations, while the dustpan #1 and the brushes #1-2 are used for the reproduction. The CT is successful if the head of the brush aligns parallel above the edge of the dustpan.

## B. Evaluation of K-VIL's Task Representation

As discussed in Section IV-B, K-VIL's task representation can be acquired from one or a few demonstration videos based

on the distance and variance criteria. Therefore, the number of demonstrations and the variations in object poses and shapes play an essential role. In this evaluation, we are interested in the following questions: (i) How do task representations learned from a different number of demonstrations affect the performance of task reproduction? In other words, what are the limitations of the task representations learned from scarce demonstrations? (ii) How many demonstrations are required to learn generalizable task representations? (iii) How do the shape and pose variations contribute to the successful extraction of such task representations? We evaluate K-VIL in one-shot and few shots imitation learning setups in Section VI-B1 and Section VI-B3, respectively. We finally answer the above questions in Section VI-B4.

**1) One-shot Imitation Learning:** We first apply K-VIL to one-shot imitation learning scenarios, where one demonstration ( $N = 1$ ) is provided for each task. As previously explained, when a single demonstration is provided, K-VIL learns a task representation based on the distance criteria of Section IV-B1, resulting in a set of 3 linear p2p constraints.

**Task extraction:** The insertion task is first learned from a single demonstration consisting of inserting the stick #4 of Fig. 10f in a paper roll. In this task, the master object is the paper roll, and the slave object is the stick. As shown in Fig. 13a, K-VIL extracts 3 keypoints subject to p2p constraints on the stick. Note that the local frame  $\mathcal{F}$  and the keypoint  $k_1$  are located near the contact point, and  $k_2$  is the farthest point on the stick from the paper roll (i.e., the master object). Similarly, in Figs. 11a, 12a, 14a and 15a, local frames are constructed on the master objects (i.e., the kettle, the tissue box, the teacup, and the rack, respectively) and 3 p2p constraints are extracted to fully constrain the pose of the slave objects (i.e., the hand, the hand, the kettle and the hat, respectively).

**Task reproduction:** As described in Section V-C, the priorities of the 3 keypoints are ranked as  $\text{Pri}_1 > \text{Pri}_2 > \text{Pri}_3$ . This respects the fact that  $k_1$  is usually the contact point of two objects, and allows the KAC to reproduce the motion of  $k_1$  more accurately than the motions of  $k_2$  and  $k_3$ . The first five columns of Table II show examples of reproduction of the insertion task IS learned from a single demonstration (task N: IS ①), as well as reproductions obtained by removing the priority from KAC as an ablation study (task O: IS<sub>np</sub> ①). For illustration purposes, we display the cases with a short stick (#10), a long stick (#6, longer than the sticks used in the demonstration), and an extra long stick (the concatenation of #8 and #9). The largest length difference is  $\sim 300$  mm. The images in the first row show the ability of K-VIL to correctly adapt the task representations of an insertion task in new scenes. This includes constructing the local frame  $\mathcal{F}$  on the master object (the paper roll) and the three p2p geometric constraints in  $\mathcal{F}$ , identifying the keypoints on the slave object (the sticks) and generating the VMP trajectories. It is important to notice that, without priority in the KAC, the keypoint  $k_1$  in IS<sub>np</sub> ①-short is not able to reach its target as accurately as in IS ①-short. Moreover, as opposed to IS<sub>np</sub> ①-long, IS ①-long is successfully executed without collision between the stick and the paper roll, thanks to the priority in KAC. However, both IS ①-ext. long and IS<sub>np</sub> ①-ext. long result in a collision

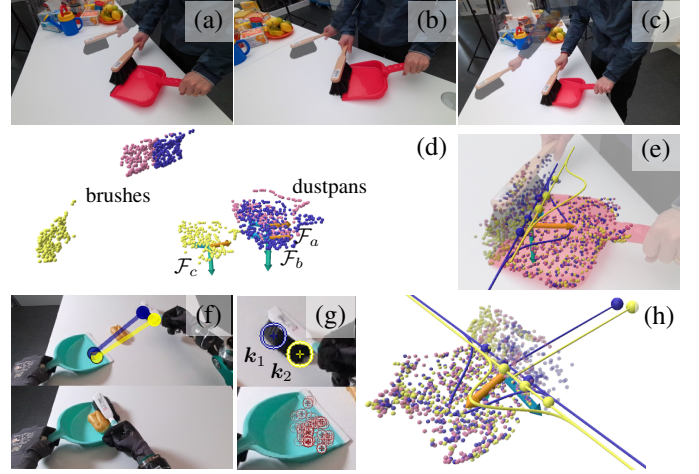


Fig. 16: K-VIL handles viewpoint mismatch in the three demonstrations (a)-(c) by aligning the corresponding local frames on the master object dustpan in (d), which results in an aligned viewpoint in (e). Two p2l constraints and their probability density functions on the principal lines are visualized. The robot reproduces the CT ③ task from a new viewpoint with a novel brush and dustpan (f), with the keypoints ( $k_1$ ,  $k_2$ ) detected on the brush hair (g). The local frame on the dustpan is determined by the  $Q = 50$  neighboring points as shown in (g). (f) and (h) depict the keypoints and their movement primitives in 2D and 3D respectively.

between the stick and the paper roll during execution due to the extended length of the stick compared to the demonstrations. The collision is more acute without priority in the KAC.

Moreover, one-shot VIL may generally fail when the learned geometric constraints are not reachable. This problem is exacerbated when demonstrations are provided from a third-person view. For example, consider Fig. 11a and Fig. 12a, that show the task representations learned from a single third-person-view demonstration of the PB and FT task, respectively. The reproductions of such task representations fail (see tasks A: PB ① and D: FT ① in Table III) due to unreachable target keypoint positions. This also indicates that task representations learned from one demonstration are not necessarily generalizable enough for motion reproduction.

Despite a few failures in the execution, K-VIL's task representations are reliably adapted to new scenes. In other words, K-VIL is able to successfully identify the keypoint positions, locate their targets, and generate the corresponding VMPs by learning their representation thanks to the combination of the proposed task representation with dense visual correspondence models. It is worth noting that this is already achieved by learning the corresponding representation from a single demonstration. Moreover, thanks to the prioritized KAC, K-VIL can handle shape variations in categorical objects via the extrapolation of the keypoint target positions. However, it cannot cope with very large shape variations. In other words, providing a single demonstration limits the learning of embodiment-independent generalizable task representations. Therefore, we then evaluate the performance of K-VIL in the case where several demonstrations are available.



Tasks	N: IS ①	O: IS <sub>np</sub> ①	N: IS ①	N: IS ①	O: IS <sub>np</sub> ①	P: IS ③	Q: IS <sub>np</sub> ③	P: IS ③	Q: IS <sub>np</sub> ③
stick	short	short	long	ext. long	ext. long	short	short	ext. long	ext. long
	3×p2p					p2p, p2l			
TR									
Reproduction									

TABLE II: Reproduction of the insertion tasks with/without priorities. Given an image of the scene before execution, K-VIL’s task representation (TR) of each task is used to identify the local frame  $\mathcal{F}$ , the keypoints ( $\bullet k_1$ ,  $\bullet k_2$ ,  $\bullet k_3$ ), their targets positions ( $\bullet k_1^g$ ,  $\bullet k_2^g$ ,  $\bullet k_3^g$ ), their movement primitives ( $\text{---}$ ,  $\text{---}$ ,  $\text{---}$ ), and the line principal manifold ( $\text{---}$ ) in tasks P and Q. The short, long and extremely long sticks correspond to sticks #10, #6, and the concatenation of #8 and #9. Task names and statistics are listed in Table VIII. The subscript np indicates that the task was reproduced without priority in KAC. The figures in each column are from one of the 20 trials for each task.

Tasks	A: PB ①	B: PB ③	C: PB ④	D: FT ①	E: FT ③	F: FT ④
	3×p2p	p2p, p2l	p2p	3×p2p	p2p, p2l	p2p
TR						
Reproduction						

TABLE III: Reproductions of tasks PB and FT learned from a third-person view without enough demonstrations lead to failure, due to unreachable geometric constraints, see tasks A, B, D, E. Generalizable task representations of PB and FT learned from enough demonstrations can be successfully executed in C and F. Legend as in Table II.

2) *Handling viewpoint mismatch*: When demonstration videos are collected from different viewpoints, e.g., in the CT ③ task in Fig. 16, we first align the demonstrations into a common viewpoint by projecting the motions of the slave objects (e.g., the brush) into each candidate local frame  $\hat{\mathcal{F}}_j$  on the master object (e.g. the dustpan) (see also Section IV-B1). In all  $|\mathcal{P}_m|$  aligned common viewpoints, we apply PCE to extract the task representations. This solves the viewpoint mismatch problem of Fig. 16d. The geometric constraints become obvious in the aligned viewpoint as shown in Fig. 16e. K-VIL extracts two p2l constraints for the CT ③ task, the combination of which forms a parallel constraint. The estimated probability density functions of the two keypoints on the corresponding principal lines ensure their target positions to be above the edge of the dustpan. It is important to

note that, not only the demonstrations can be recorded from different viewpoints, but also the reproduction of the learned skill by the robot can be performed from a viewpoint that is significantly different from any demonstration, as shown in Figs. 16f and 16h. For the seek of clarity, we discuss the results of K-VIL in the aligned viewpoints in the remaining evaluations, although the demonstrations and reproductions happen in different viewpoints as discussed for the CT ③ task.

3) *Updating Constraints Incrementally*: Here, we apply K-VIL to few-shot imitation learning scenarios where additional demonstrations are incrementally provided for each task. In this case, K-VIL is trained based on the variance criteria to learn more generalizable task representations based on various linear and non-linear constraints (see Sections IV-B2- IV-B3). Figs. 11 to 15 show the task representations of each task in



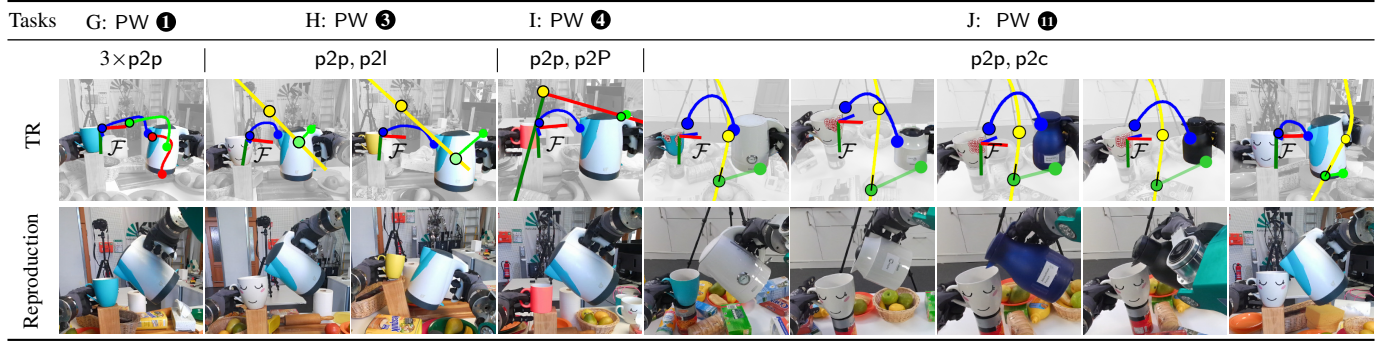


TABLE IV: Reproductions of tasks PW learned from different number of demonstrations. For task H, I and J, we mark the learned principal manifold  $\mathcal{M}_{\text{line}}$  (—),  $\mathcal{M}_{\text{plane}}$  (—, —) and  $\mathcal{M}_{\text{curve}}$  (—), respectively. Additionally, the point (●) indicates the mean of the demonstrated targets of keypoint  $k_2$  on the corresponding principal manifolds. Other legends as in Table II.

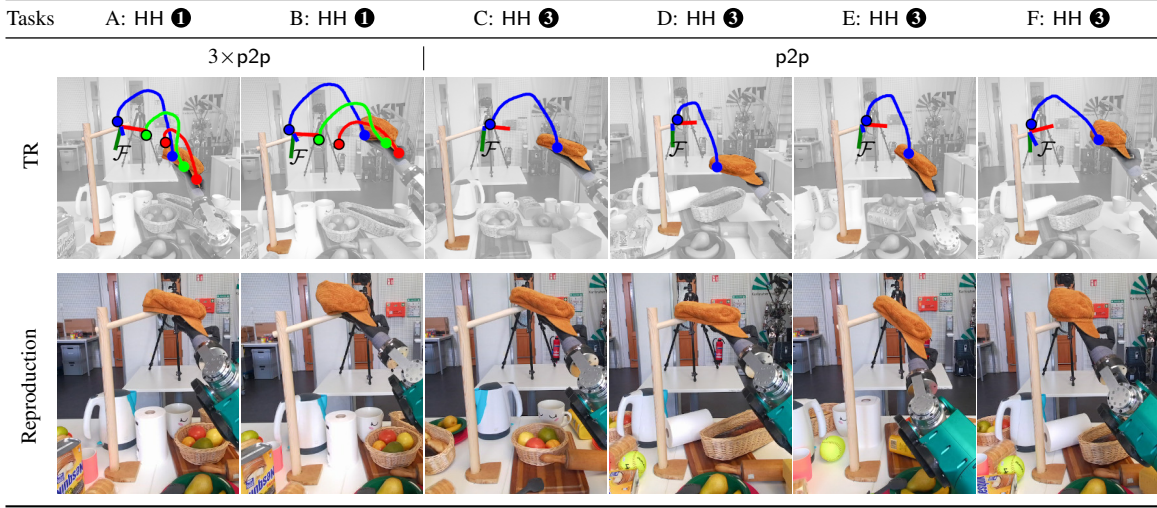


TABLE V: Reproductions of tasks HH. Legend as in Table II.

$\mathcal{T}_E$  learned by K-VIL from several demonstrations.

#### Task extraction and reproduction of insertion tasks IS.

Providing several demonstrations allows us to consider variations in the demonstrated task, and thus to extract prioritized geometric constraints. For example, the keypoint  $k_1$  in Fig. 13b is the most invariant point on the stick across all demonstrations, while  $k_2$  is subject to a p2l constraint. Note that this contrasts with Fig. 13a, where all keypoints were subject to p2p constraints. With such task representation, KAC successfully handles all stick lengths by fulfilling the p2p and p2l constraints, as shown in Table II for IS ❸ (P). Despite a drop in accuracy and precision (see Section VI-C, Table VIII), the KAC without priority still leads to successful task completion in this case (see IS<sub>np</sub> ❸ (Q) in Table II). Overall, the extrapolation abilities of K-VIL are significantly increased by providing 3 demonstrations instead of 1. Note that, due to the nature of the p2l constraint and the priority mechanism in KAC, sticks of arbitrary length can be handled. As the line manifold goes through both  $k_1$  and  $k_2$ , K-VIL implicitly learns a colinear constraint for the two keypoints.

#### Task extraction and reproduction of PB and FT tasks.

For these two tasks, 3 demonstrations are not sufficient to completely represent the task, and K-VIL may coincidentally extract a superfluous p2l constraint (see  $\mathcal{M}_{\text{line}}$  in Fig. 11b and

Fig. 12b). This forces the robot to place its hand similarly as demonstrated by the human. However, due to the third-person view adopted for the demonstration, this cannot be achieved by the robot, therefore resulting in failed executions of the tasks PB ❸ and FT ❸ as for PB ❶ and FT ❶ (see Table III). The unnecessary p2l constraints are removed by providing K-VIL with an additional demonstration (see Figs. 11c and 12c), allowing the tasks to be successfully reproduced (see PB ❹ and FT ❹ in Table III).

**Task extraction and reproduction of PW tasks.** Similarly to PB and FT tasks, K-VIL's representation obtained from 3 demonstrations for PW results in superfluous p2p and p2l constraints. Although the task may still be executed by the robot (see PW ❸ in the second column of Table IV), the superfluous constraints may lead to collisions between the kettle and the environment in other cases. For example, in the third column in Table IV), with some specific initial poses of the kettle, the generated VMPs and the line constraints lead to rotation of the kettle in reversed direction during the execution, thus causing reproduction failures. These restrictive task representations are alleviated by providing K-VIL with an additional demonstration. By doing so, the problematic constraints are updated to different types, e.g., the p2l constraint in Fig. 14b becomes a p2P constraint in Fig. 14c.

This significantly improves K-VIL’s extrapolation abilities, as the p2P constraint is necessary to constrain the kettle in a vertical plane while being less restrictive than the previous p2l constraint. Notice that the density force within the plane constraint ensures the tilting angle of the kettle is similar to the demonstrations. A successful reproduction of such task representations is shown in Table IV (PW ④). When enough demonstrations are provided ( $N = 11$  in Fig. 14d), K-VIL instead extracts a p2c constraint for the keypoint  $k_2$  at the bottom of the kettle. The reproduction results with different types of kettles are shown in the last five columns of Table IV (PW ⑤). Intuitively speaking, the p2c constraint aligns better with our understanding of a pouring task. Moreover, the last column in Table IV shows that, given a tilted initial pose of the kettle, KAC is able to correct the pose of the kettle in the end. The set of PW tasks in Table IV also demonstrates that the proposed approach generalizes well to categorical objects with different colors, sizes and shapes, and is robust to change of background and viewpoints.

**Task extraction and reproduction of HH tasks.** Unlike PB, FT and PW tasks, K-VIL’s representation obtained from 3 demonstrations for HH does not result in superfluous p2l constraints. Instead, due to the obvious pose variations in the hat, K-VIL consistently extracts a single keypoint  $k_1$  on the backside of the hat with a p2p constraint, which encodes the demonstrated position invariances observed in the local frame near the end of the hanging stick. As shown in Table V, the target poses of the hat in HH ⑤ vary according to its different initial poses, while in HH ① they are fully determined by the 3 p2p constraints. Figs. 15c to 15e show the influence of the shape variations on the selection of local frames  $\mathcal{F}$ .

4) *Evaluation summary:* As shown by our experiments, K-VIL’s task representations allow the successful learning of diverse tasks and their reproduction in new cluttered scenes with large shape and pose variations in categorical objects. As opposed to [23], [24], our approach is not constrained to conserving the same viewpoint between demonstrations and reproductions, and thus is more flexible. Here, we further discuss the influence of the number and diversity of the demonstrations on such task representations.

**Limitations of scarce demonstrations.** The task representations learned from scarce demonstrations hinder the performance and the extrapolation ability of K-VIL in three ways, namely, (i) they may be embodiment-dependent, and thus cannot be reproduced by the robot, e.g., in Table III for PB ①, PB ③, FT ① and FT ③; (ii) they may lead to collisions during the execution, e.g., in IS ③-ext. long (see Table II) and PW ③ for improper kettle start pose (see Table IV); (iii) when reproduced successfully, the control accuracy is reduced compared to the task representations learned from more demonstrations (see Section VI-C for details). These limitations motivate us to evaluate the number of demonstrations that are required to learn a generalizable task representation according to the criteria of Section VI-A for each considered task.

**Adequate number of demonstrations.** Table VI summarizes the extracted geometric constraints of the five extraction tasks for different number  $N$  of demonstrations. It is interesting to notice that the learned task representations do

		Number of demonstrations (N)				
$\mathcal{T}_E$		1	3	4	5	11
a	PB	$3 \times \text{p2p}$	p2p, p2l	p2p	p2p	p2p
b	FT	$3 \times \text{p2p}$	p2p, p2l	p2p	p2p	p2p
c	PW	$3 \times \text{p2p}$	p2p, p2l	p2p, p2P	p2p, p2P	p2p, p2c
d	HH	$3 \times \text{p2p}$	p2p	p2p	p2p	p2p
e	IS	$3 \times \text{p2p}$	p2p, p2l	p2p, p2l	p2p, p2l	p2p, p2l
f	CT	$3 \times \text{p2p}$	p2l, p2l	p2l, p2l	p2l, p2l	p2l, p2l

TABLE VI: Extraction tasks and the geometric constraints of each task learned from different number of demonstrations. We mark the cases (   ) where the learned task representations converge, and highlight the cases (in blue) where the learned task representations are generalizable.

not change anymore, i.e., converge, after a certain number of demonstrations, e.g.,  $N = 4$  for PB and FT,  $N = 11$  for PW, and  $N = 3$  for HH and IS. Moreover, the task representations become generalizable almost at the same time as they converge. As an exception, the representations of PW already generalize with a p2p and a p2P constraint learned from 4 demonstrations. Importantly, the learned geometric constraints do not have to include a p2p constraint, see e.g., the CT task that is represented by two p2l constraints. Overall, this demonstrates that our approach efficiently extracts generalizable task representations from considerably less demonstrations than state-of-the-art approaches.

**Object pose and shape variations.** Importantly, K-VIL’s task representations rely on the variations observed in the demonstrations to extract appropriate constraints. For example, in PB, the pose variations of the demonstrator’s hand (Fig. 11c) allow K-VIL to distinguish the tip of the middle finger (used to press the button) from the other candidate points on the hand. In other words, these variations enable the efficient extraction of the keypoint  $k_1$  subject to a p2p constraint. The spatial distribution of keypoints across demonstrations is also obtained via pose variations. For instance, this allows K-VIL to associate the keypoint  $k_2$  with geometric constraints such as p2l (Figs. 14b and 16e), p2P (Fig. 14c) and p2c (Fig. 14d).

Shape variations also facilitate the efficient extraction of keypoints and geometric constraints. For example, the stick length variations in the insertion task IS allow K-VIL to extract the keypoints  $k_1$  and  $k_2$  subject to a p2p and p2l constraint, respectively, as intuitively shown in Fig. 13b. In other cases, shape variations in the master objects help to remove redundancy in canonical local frames. For example, all the  $J = 300$  canonical local frames on the rack in task HH of Fig. 15b are equivalent. This is due to the absence of variations in the rack across the provided demonstrations. In this case, K-VIL selects the local frame  $\mathcal{F}$  as the closest on average to the keypoint  $k_1$  on the hat. This redundancy is removed by introducing shape variations in the master objects. For example, by considering two variations of the rack in the task HH as in Fig. 15c,  $k_1$  is position-invariant only if it is represented in the local frames near the contact point between the rack and the hat. This allows K-VIL to focus on these local frames and to filter out the others. We showed in Fig. 15d that K-VIL’s representation converges and remains the same

index	$\mathcal{T}_E$	role	object	PV	SV	TR
1	PB	master slave	kettle hand	- ✓	✗ ✗	Figs. 11b and 11c
2	FT	master slave	tissue box hand	- ✓	✓ ✗	Figs. 12b and 12c
3	PW	master slave	teacups kettle	- ✓	✓ ✓	Figs. 14b to 14d
4	HH	master slave	rack hat	- ✓	✗ ✗	Fig. 15b
5	HH	master slave	rack hat	- ✓	✓ ✗	Fig. 15c
6	HH	master slave	rack hat	- ✓	✗ ✓	Figs. 15d and 15e
7	IS	master slave	paper roll stick	- ✗	✗ ✓	Figs. 13b and 13c
8	CT	master slave	dustpan brush	- ✓	✓ ✓	Fig. 16

TABLE VII: Pose variations (PV) and shape variations (SV) in the demonstrations along with the detected master and slave objects. The corresponding task representations (TR) are linked in the last column. Pose variations of the master objects are not relevant (-) as the local frames representing the object pose are constructed on the masters.

as Figs. 15b and 15c even if more demonstrations with large shape and pose variations in the hats are available. Similarly, shape variations in the tissue boxes in FT allow the selection of local frames around the grasping point, while shape variations in the teacups (PW) ensure that local frames for the pouring task are around their rim. A summary of the effect of pose and shape variations in the considered tasks is given by Table VII. Overall, pose or/and shape variations of master and slave objects in the demonstrations are not only handled by K-VIL, but also facilitate the efficient, joint extraction of local frames, keypoints, and geometric constraints.

### C. Evaluation of KAC

In Section VI-B, we demonstrated K-VIL's ability to extract generalizable task representations from a small number of demonstrations. In particular, we showed that these task representations successfully adapt to new cluttered scenes with categorical objects, regardless of whether the task can be reproduced by the robot. Therefore, in this section, we evaluate the proposed KAC in terms of control accuracy, control precision (i.e., repeatability) and success rate on the 18 tasks described in Section VI-A (see also Table VIII). Each task is reproduced  $N_r = 20$  times according to the evaluation protocol in Section VI-A. For each task, we record the trajectories of all relevant keypoints during the execution period. Since we are particularly interested in the regulation behavior of KAC when the keypoints satisfy their corresponding geometric constraints, i.e., when the VMPs finish, we record the keypoint trajectories for an additional 2 s time window, denoted by  $\mathcal{T}_{\text{end}}$ . These trajectories are then compared to the corresponding keypoints' target trajectories, which correspond to the VMP trajectories for keypoints subject to p2p constraint, and to the attractor trajectories for keypoints subject to other types of constraints. Notice that, for the latter, the 3D position of the attractor is recovered from the corresponding 1-dimensional

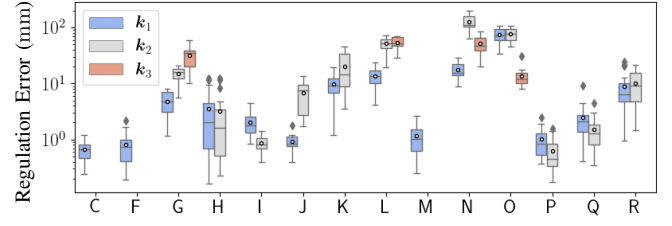


Fig. 17: Regulation errors between each keypoint and its target estimated in  $\mathcal{T}_{\text{end}}$  over  $N_r = 20$  trials for tasks C to R in Table VIII. The mean regulation error, i.e., control accuracy Acc., is depicted as  $\circ$ . The box shows the first and third quartiles of the regulation error of each keypoint, with the bar inside it indicating the median.

VMP in the orthogonal direction of the corresponding principal manifold (see also Section V-A).

First, we evaluate the ability of KAC to satisfy the learned keypoints' geometric constraints. To do so, we compute the regulation error of each keypoint during  $\mathcal{T}_{\text{end}}$  of each trial as

$$e_v = \frac{1}{T_r} \sum_{t=0}^{T_r} \|\mathbf{k}_{l,v}^g(t) - \mathbf{k}_{l,v}(t)\|_2, \quad v \in [1, N_r],$$

where  $T_r$  is the total timesteps recorded in  $\mathcal{T}_{\text{end}}$ , and  $\mathbf{k}_{l,v}, \mathbf{k}_{l,v}^g$  are the recorded and target positions of the considered keypoint in the  $v^{\text{th}}$  trial, respectively. Fig. 17 displays the distribution of the keypoint's regulation errors for 20 trials of each task, where the mean values ( $\circ$ ) correspond to the control accuracy

$$\text{Acc.} = \frac{1}{N_r} \sum_{v=0}^{N_r} e_v.$$

Second, we evaluate the control precision (i.e., repeatability) of KAC for all keypoints and all tasks. It is computed as

$$\text{Prec.} = \sqrt{\frac{1}{N_r} \frac{1}{T_r} \sum_{v=1}^{N_r} \sum_{t=1}^{T_r} \|\mathbf{k}_{l,v}(t) - \boldsymbol{\mu}_{\mathbf{k}_{l,v}}\|_2^2}, \quad (5)$$

where we defined  $\boldsymbol{\mu}_{\mathbf{k}_{l,v}} = \frac{1}{T_r} \sum_{t=1}^{T_r} \mathbf{k}_{l,v}(t)$ . Finally, we also report the success rate obtained for each task according to the evaluation protocols described in Section VI-A.

Table VIII presents the evaluation results of KAC in terms of the three aforementioned metrics. As suggested by the qualitative evaluations in Table III, when learned from fewer than 4 demonstrations, the tasks PB and FT result in a 0% success rate (see A-B, D-E in Table VIII). As discussed in Section VI-B3, this is due to geometric constraints that are unreachable for the robot. In contrast, the task representations learned from 4 demonstrations are generalizable (see Table VI in Section VI-B4), and thus KAC reaches sub-millimeter control accuracy and precision, as well as  $\geq 90\%$  success rates (see C and F in Table VIII).

As shown in Fig. 17 and Table VIII (G-J), PW ④ and PW ⑪ outperform PW ③ in terms of control accuracy, precision, and success rate. This is due to the fact that, in contrast to PW ③, PW ④ and PW ⑪ are generalizable (see Table VI and Table IV). Although PW ① displays a relatively high control precision and success rate, its control accuracy remains low and the pose of the kettle is fully constrained by 3 p2p constraints, which hinders K-VIL's extrapolation abilities. Interestingly, while  $\mathbf{k}_1$ 's control accuracy in the task



PW increases with the number of demonstrations,  $k_2$ 's highest control accuracy is obtained in PW ④ (I). This is due to  $k_2$ 's p2P constraint in PW ④, which is easier to fulfill than the p2c constraint in PW ①. We observe a lower precision in PW ④ than in PW ① for the same reason.

Similarly as in PW, HH ③ (M) and IS ③ (P) are reproduced with higher control accuracy and precision than HH ① (L) and IS ① (N), respectively. Moreover, thanks to the priority introduced in KAC, HH ① (L) and IS ① (N) achieve 95% and 62% success rate, respectively, despite the single available demonstration. The lower performance of IS ① is explained by the fact that we consider an extremely long stick, which requires high generalization capabilities, in the reproduction (see Table II). In contrast, the hat in HH is only slightly deformable and thus results in low shape variations. Notice that, for HH and IS, we do not compare the control accuracy between keypoints subject to different geometric constraints (in gray in Table VIII). As these keypoints are controlled according to different priorities within KAC, their reported accuracy highly depends on the shape variations occurring in these two tasks. For example, when using sticks of various lengths in the insertion task IS ① (N), KAC assigns the highest priority to the keypoint  $k_1$ , so that  $k_1$  is obviously controlled with higher accuracy than  $k_2$  and  $k_3$ . Interestingly,  $k_2$ 's highest regulation error in this task is  $\sim 150$  mm, which corresponds to the maximum length difference between the sticks used in the demonstration and in the reproduction. Therefore, in HH and IS, the control accuracy of  $k_2$  and  $k_3$  is rather affected by the experiment setups (e.g., stick lengths) than by KAC. Since the two p2l constraints in CT ③ share the same priority, the two keypoints are equally controlled towards the region with high likelihood on the corresponding principal lines. The two spring-damper systems for the two keypoints result in equilibrium. Therefore, the control accuracy and the precision of the two keypoints are comparable.

Compared to IS ①, the ablation study IS<sub>np</sub> ① (O), conducted without KAC's priorities, shows a drop in success rate (25%), as  $k_3$  achieves the highest accuracy at the expense of  $k_1$  and  $k_2$ . This is expected as the three identical virtual spring-damper systems for  $k_1$ ,  $k_2$  and  $k_3$  are in equilibrium, and  $k_3$  usually locates in the middle of  $k_1$  and  $k_2$ . As the ablation tasks PW<sub>np</sub> ① (K) and IS<sub>np</sub> ③ (Q) are reproduced from generalizable task representations, removing KAC's priorities does not affect their success rate. However, the control accuracy and precision drop compared to PW ① (J) and IS ③ (P).

## VII. DISCUSSION

In this paper, we proposed the novel keypoints-based visual imitation learning (K-VIL) approach that learns sparse, object-centric, and embodiment-independent task representations from a small set of demonstration videos. K-VIL's task representations are based on the extraction of geometric constraints by a PCE, which covers a wide range of constraints. The proposed PCE enables one-shot and few-shot VIL and updates the learned task representations when additional demonstrations are incrementally provided, thus endowing them with enhanced extrapolation capabilities. K-VIL's task representations also include task-specific keypoint

	$\mathcal{T}_R$	TR	Acc. (mm)			Prec. (mm)			R (%)
			$k_1$	$k_2$	$k_3$	$k_1$	$k_2$	$k_3$	
A	PB ①	Fig. 11a	×	×	×	×	×	×	0
B	PB ③	Fig. 11b	×	×	×	×	×	×	0
C	PB ④	Fig. 11c	0.67	-	-	0.34	-	-	90
D	FT ①	Fig. 12a	×	×	×	×	×	×	0
E	FT ③	Fig. 12b	×	×	×	×	×	×	0
F	FT ④	Fig. 12c	0.82	-	-	0.65	-	-	100
G	PW ①	Fig. 14a	4.81	15.01	32.09	0.75	0.78	0.72	95
H	PW ③	Fig. 14b	3.56	3.21	-	1.20	2.12	-	87
I	PW ④	Fig. 14c	2.01	0.88	-	1.04	5.67	-	94
J	PW ①	Fig. 14d	0.93	6.80	-	0.42	0.56	-	100
K	PW <sub>np</sub> ①	Fig. 14d	9.71	20.08	-	0.90	1.16	-	100
L	HH ①	Fig. 15a	13.53	51.33	54.28	0.67	0.67	0.67	95
M	HH ③	Fig. 15b	1.48	-	-	0.47	-	-	95
N	IS ①	Fig. 13a	17.77	125.40	51.98	1.55	1.30	1.34	62
O	IS <sub>np</sub> ①	Fig. 13a	74.55	77.05	13.40	1.03	0.97	1.02	25
P	IS ③	Fig. 13b	1.01	0.63	-	0.36	0.36	-	95
Q	IS <sub>np</sub> ③	Fig. 13b	2.46	1.50	-	2.32	1.39	-	90
R	CT ③	Fig. 16e	8.89	10.22	-	2.04	2.03	-	95

TABLE VIII: Evaluation of KAC in terms of control accuracy (Acc.), precision (Prec.), and success rate (R) for each category of tasks with task representations (TR) learned from different numbers  $N$  of demonstrations. Ablation studies (denoted by <sub>np</sub>) are also conducted in tasks K, O, and Q by removing the priority (see Section V-C) from KAC. The cases where data is not available due to failure execution and where a specific keypoint is not required are denoted as  $\times$  and  $-$ , respectively. Gray numbers in N and O indicate inconsequential values.

control policies encoded as VMPs, which are leveraged for task execution by a prioritized keypoint-based admittance controller (KAC). Compared to control policies based on RL or on visual servoing, VMPs allow a flexible temporal scaling and support via-points (including start and target position) adaptation. Therefore, they crucially contribute to K-VIL's generalization capabilities by extrapolating the keypoint target positions on the learned principal manifold.

As highlighted in our evaluation, K-VIL consistently learned generalizable task representations for six daily manipulation tasks, which involved highly cluttered scenes, new instances of categorical objects, and large variations in object poses and shapes. Importantly, we showed that the learned task representations converges and becomes generalizable with significantly fewer demonstrations than state-of-the-art approaches such as [12], [14], [22], [27]. Interestingly, the sparse keypoint-based geometric constraints extracted by K-VIL mostly aligned with human intuition. This includes the extraction of a single p2p constraint for pressing a button, of a pair of p2p and p2l constraints for the insertion task, and of a p2p coupled with a p2c constraint for the pouring task, among others.

It is important to emphasize that the decomposed control and priority mechanism of the KAC allowed us to endow K-VIL with reliable extrapolation capabilities. Indeed, our quantitative evaluations demonstrated K-VIL's ability to reproduce the learned task representations with high control accuracy, control precision, and success rate. Particularly,

K-VIL accurately handled very large shape variations in the considered insertion task. In contrast, previous works did not or only briefly discuss the extrapolation capabilities of their approaches [12], [16], [23], [19]. For instance, Jin & Jagersand [12] only showed extrapolation to another instance of the hammer category with very small shape variation without providing any quantitative evaluations.

It is important to note that the variations in object poses and shapes play an essential role in learning generalizable task representations. This is even more relevant when only a small number of demonstrations are provided. Without such variations, K-VIL can still generalize to categorical objects thanks to the dense visual descriptors, but achieve lower control accuracy, precision, and success rate, and may fail in some extreme cases, e.g., in the one-shot VIL setup.

#### A. Limitations and Future Work

K-VIL imposes limitations in terms of visual perception models and task representations. On the one hand, we assume that the keypoints are on the surface of objects and omit transparent, reflective, and thin objects (note that this is also discussed in [20], [42]). This hinders K-VIL from being used in many real-world tasks. Furthermore, all keypoints must be visible in the demonstrations, which may not always be enforced in reality. In other words, K-VIL learns from demonstrations with and without viewpoint mismatch, as long as the keypoints of interest are not occluded. In the long run, we believe that the dense correspondence models should be combined with state-of-the-art scene representation models (e.g., [42]) or with point generative models (e.g., [43]) for better correspondence detection and for tackling the occlusion problems. This would allow the imitator to observe objects that are visually more challenging and to learn the task from demonstrations with (self-)occlusions.

It is worth noticing that K-VIL's keypoints correspond to the sub-symbolic parameters of a motion. Therefore, they do not necessarily have a clear semantic interpretation, which is also important for learning comprehensive task models. Bridging the gap between the symbolic and sub-symbolic levels remains an important challenge in (visual) IL. Importantly, the symbolic representation of a task [44], [45] also has limitations, which can be alleviated by integrating sub-symbolic information. For example, a contain affordance in a pouring task implies that the opening of the spout of the kettle should be placed above the contain affordance region [46]. However, this semantic representation alone cannot describe different types of pouring: For example, pouring beer requires tilting the glass and aligning the beer with the side of a glass. Instead, additional sub-symbolic parameters would allow realizing specific styles of task execution.

In this sense, K-VIL deals with the sub-symbolic part of the task. Namely, its ability to update the geometric constraints allows us (i) to reproduce a task with a specific style, and (ii) to eliminate unnecessary keypoints and geometric constraints and to update the distribution of the keypoints on the extracted constraints when more demonstration styles are available. K-VIL may then be augmented with an extraction method [47]

to estimate the links between the extracted keypoints and the symbolic task representation. For instance, the probability distribution of the keypoints on their principal manifolds may be used to determine the affordance regions [46], [48], the spatial relations [49], and the grasping or effect points [10], [47]. We will investigate these aspects in our future work.

In this paper, we only considered uni-manual manipulation tasks that can be modeled as the combination of five basic geometric constraints in Fig. 4 in a single layer of master-slave relationship. As future work, we plan to extend K-VIL to include other types geometric constraint and to bimanual manipulation tasks by considering bimanual coordination strategies [50] and a hierarchy of master-slave relationships. Moreover, we will extend K-VIL for periodic motions such as stirring or wiping motions [23], as well as for handling articulated objects [51], [52].

#### REFERENCES

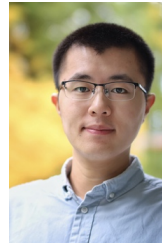
- [1] A. Bandura and R. H. Walters, *Social learning theory*. Englewood cliffs Prentice Hall, 1977, vol. 1.
- [2] D. M. Wolpert, J. Diedrichsen, and J. R. Flanagan, "Principles of sensorimotor learning," *Nature Reviews Neuroscience*, pp. 739–751, 2011.
- [3] C. J. Burke, P. N. Tobler, M. Baddeley, and W. Schultz, "Neural mechanisms of observational learning," *Proceedings of the National Academy of Sciences*, pp. 14 431–14 436, 2010.
- [4] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *ICCV*, 2017, pp. 2980–2988.
- [5] P. Florence, L. Manuelli, and R. Tedrake, "Dense Object Nets: Learning dense visual object descriptors by and for robotic manipulation," in *CoRL*, ser. Proceedings of Machine Learning Research, vol. 87, 2018, pp. 373–385.
- [6] K. Meng and A. Eloyan, "Principal manifold estimation via model complexity selection," *Journal of the Royal Statistical Society. Series B. Statistical methodology*, vol. 83, no. 2, pp. 369–394, 2021.
- [7] W. Gao and R. Tedrake, "kPAM 2.0: Feedback control for category-level robotic manipulation," *IEEE RA-L*, vol. 6, no. 2, pp. 2962–2969, 2021.
- [8] R. Xu, F.-J. Chu, C. Tang, W. Liu, and P. Vela, "An affordance keypoint detection network for robot manipulation," *IEEE RA-L*, vol. 6, no. 2, pp. 2870–2877, 2021.
- [9] M. Sharma and O. Kroemer, "Generalizing object-centric task-axes controllers using keypoints," in *IEEE ICRA*, 2021, pp. 7548–7554.
- [10] Z. Qin, K. Fang, Y. Zhu, L. Fei-Fei, and S. Savarese, "KETO: Learning keypoint representations for tool manipulation," in *IEEE ICRA*, 2020, pp. 7278–7285.
- [11] M. Sharma, J. Liang, J. Zhao, A. Lagrassa, and O. Kroemer, "Learning to compose hierarchical object-centric controllers for robotic manipulation," in *CoRL*, vol. 155. PMLR, 2021, pp. 822–844.
- [12] J. Jin and M. Jagersand, "Generalizable task representation learning from human demonstration videos: a geometric approach," in *IEEE ICRA*, 2022, pp. 2504–2510.
- [13] Y. Zhou, J. Gao, and T. Asfour, "Learning via-point movement primitives with inter- and extrapolation capabilities," in *IEEE/RSJ IROS*, 2019, pp. 4301–4308.
- [14] P. Sharma, D. Pathak, and A. K. Gupta, "Third-person visual imitation learning via decoupled hierarchical controller," in *NeurIPS*, 2019.
- [15] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine, "AVID: Learning multi-stage tasks via pixel-level translation of human videos," in *R:SS*, 2020.
- [16] M. Sieb, Z. Xian, A. Huang, O. Kroemer, and K. Fragkiadaki, "Graph-structured visual imitation," in *CoRL*, vol. 100. PMLR, 2020, pp. 979–989.
- [17] Y. Liu, A. Gupta, P. Abbeel, and S. Levine, "Imitation from observation: Learning to imitate behaviors from raw video via context translation," in *IEEE ICRA*, 2018, pp. 1118–1125.
- [18] D. Dwibedi, J. Tompson, C. Lynch, and P. Sermanet, "Learning actionable representations from visual observations," in *IEEE/RSJ IROS*, 2018, pp. 1577–1584.
- [19] H. Karman, F. Torabi, G. Warnell, and P. Stone, "Adversarial imitation learning from video using a state observer," in *IEEE ICRA*, 2022.



- [20] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitzmann, "Neural Descriptor Fields: SE(3)-equivariant object representations for manipulation," in *IEEE ICRA*, 2022, pp. 6394–6400.
- [21] P. Florence, L. Manuelli, and R. Tedrake, "Self-supervised correspondence in visuomotor policy learning," *IEEE RA-L*, vol. 5, pp. 492–499, 2020.
- [22] D. Pathak, P. Mahmoudieh, G. Luo, P. Agrawal, D. Chen, Y. Shentu, E. Shelhamer, J. Malik, A. A. Efros, and T. Darrell, "Zero-shot visual imitation," in *ICLR*, 2018.
- [23] J. Yang, J. Zhang, C. Settle, A. Rai, R. Antonova, and J. Bohg, "Learning periodic tasks from human demonstrations," in *IEEE ICRA*, 2022, pp. 8658–8665.
- [24] J. Pari, N. M. M. Shafiullah, S. P. Arunachalam, and L. Pinto, "The surprising effectiveness of representation learning for visual imitation," in *R:SS*, 2022.
- [25] F. Torabi, G. Warnell, and P. Stone, "Imitation Learning from Video by Leveraging Proprioception," in *IJCAI*, 2019, pp. 3585–3591.
- [26] —, "Generative Adversarial Imitation from Observation," in *Proceedings of the ICML Workshop on Imitation, Intent, and Interaction*, 2019.
- [27] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine, "Time-Contrastive Networks: Self-Supervised Learning from Video," in *IEEE ICRA*, 2018, pp. 1134–1141.
- [28] H. Karnan, G. Warnell, X. Xiao, and P. Stone, "Voila: Visual-observation-only imitation learning for autonomous navigation," in *IEEE ICRA*, 2022, pp. 2497–2503.
- [29] M. Muhlig, M. Gienger, J. J. Steil, and C. Goerick, "Automatic selection of task spaces for imitation learning," in *IEEE/RSJ IROS*, 2009, pp. 4996–5002.
- [30] A. L. P. Ureche, K. Umezawa, Y. Nakamura, and A. Billard, "Task parameterization using continuous constraints extracted from human demonstrations," *IEEE T-RO*, pp. 1458–1471, 2015.
- [31] Z. Dodds, M. Jägersand, G. Hager, and K. Toyama, "A hierarchical vision architecture for robotic manipulation tasks," in *Computer Vision Systems*. Springer Berlin Heidelberg, 1999, pp. 312–330.
- [32] J. P. Hespanha, Z. Dodds, G. D. Hager, and A. S. Morse, "What tasks can be performed with an uncalibrated stereo vision system?" *IJCV*, vol. 35, no. 1, pp. 65–85, 1999.
- [33] M. Gridseth, O. Ramirez, C. P. Quintero, and M. Jagersand, "ViTa: Visual task specification interface for manipulation with uncalibrated visual servoing," in *IEEE ICRA*, 2016, pp. 3434–3440.
- [34] J. Jin, L. Petrich, M. Dehghan, and M. Jägersand, "A geometric perspective on visual imitation learning," in *IEEE/RSJ IROS*, 2020, pp. 5194–5200.
- [35] L. Manuelli, Y. Li, P. R. Florence, and R. Tedrake, "Keypoints into the future: Self-supervised correspondence in model-based reinforcement learning," in *CoRL*, 2020, pp. 693–710.
- [36] C. Lugaresi, J. Tang, H. Nash, C. McClanahan, E. Ubowaja, M. Hays, F. Zhang, C.-L. Chang, M. Yong, J. Lee, W.-T. Chang, W. Hua, M. Georg, and M. Grundmann, "MediaPipe: A framework for building perception pipelines," in *Third Workshop on Computer Vision for AR/VR at IEEE Computer Vision and Pattern Recognition (CVPR) 2019*, 2019.
- [37] A. Paraschos, C. Daniel, J. Peters, and G. Neumann, "Using probabilistic movement primitives in robotics," *Auton. Robot.*, vol. 42, no. 3, pp. 529–551, 2018.
- [38] N. Halko, P.-G. Martinsson, and J. A. Tropp, "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions," *SIAM Review*, vol. 53, no. 2, pp. 217–288, 2011.
- [39] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *JMLR*, vol. 12, pp. 2825–2830, 2011.
- [40] Q.-Y. Zhou, J. Park, and V. Koltun, "Open3D: A modern library for 3D data processing," *arXiv:1801.09847*, 2018.
- [41] T. Asfour, M. Wächter, L. Kaul, S. Rader, P. Weiner, S. Ottenhaus, R. Grimm, Y. Zhou, M. Grotz, and F. Paus, "ARMAR-6: A high-performance humanoid for human-robot collaboration in real world scenarios," *IEEE RAM*, vol. 26, no. 4, pp. 108–121, 2019.
- [42] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola, "NeRF-Supervision: Learning dense object descriptors from neural radiance fields," in *IEEE ICRA*, 2022, pp. 6496–6503.
- [43] J. Lei and K. Daniilidis, "Cadex: Learning canonical deformation coordinate space for dynamic surface representation via neural homeomorphism," in *CVPR*, 2022, pp. 6624–6634.
- [44] C. R. G. Dreher and T. Asfour, "Learning temporal task models from human bimanual demonstrations," in *IEEE/RSJ IROS*, 2022.
- [45] M. Hassanin, S. Khan, and M. Tahtali, "Visual affordance and function understanding," *ACM Computing Surveys (CSUR)*, vol. 54, pp. 1 – 35, 2018.
- [46] D. Hadjivelichkov, S. Zwane, L. Agapito, M. P. Deisenroth, and D. Kanoulas, "One-Shot Transfer of Affordance Regions? AffCorrs!" in *CoRL*, 2022.
- [47] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu, "Synergies Between Affordance and Geometry: 6-DoF Grasp Detection via Implicit Representations," in *R:SS*, 2021.
- [48] T.-T. Do, A. Nguyen, and I. Reid, "AffordanceNet: An End-to-End Deep Learning Approach for Object Affordance Detection," in *IEEE ICRA*, 2018, pp. 5882–5889.
- [49] R. Kartmann, D. Liu, and T. Asfour, "Semantic scene manipulation based on 3d spatial object relations and language instructions," in *IEEE/RAS Humanoids*, 2021, pp. 306–313.
- [50] F. Krebs and T. Asfour, "A bimanual manipulation taxonomy," *IEEE RA-L*, vol. 7, no. 4, pp. 11 031–11 038, 2022.
- [51] Z. Xu, Z. He, and S. Song, "Universal Manipulation Policy Network for Articulated Objects," *IEEE RA-L*, vol. 7, pp. 2447–2454, 2022.
- [52] R. Wu, Y. Zhao, K. Mo, Z. Guo, Y. Wang, T. Wu, Q. Fan, X. Chen, L. Guibas, and H. Dong, "VAT-mart: Learning Visual Action Trajectory Proposals for Manipulating 3D ARTiculated Objects," in *ICLR*, 2022.



**Jianfeng Gao** received the bachelor's degree in Mechatronics from Beijing Institute of Technology, China and the M.S. degree in Mechatronics and Information Technology from Karlsruhe Institute of Technology (KIT), Germany, in 2015 and 2018 respectively. Currently, he is Ph.D. student at the High Performance Humanoid Technologies Lab (H<sup>2</sup>T) at KIT. His research focuses on visual imitation learning and compliance control for bimanual manipulation.



**Zhi Tao** Zhi Tao received a bachelor's in mechatronics from the Beijing Institute of Technology, China, and an M.Sc. degree in mechatronics and information technology from the Karlsruhe Institute of Technology (KIT), Germany, in 2016 and 2021 respectively. Currently, he is Robotics Algorithm Engineer with Dreame Technology, Suzhou, China. His research interests are the recognition and use of everyday tools in daily life scenes and visual imitation learning.



**Noémie Jaquier** received the B.Sc. degree on Microengineering, the M.Sc. degree in Robotics and Autonomous Systems, and the PhD degree from the Ecole Polytechnique Fédérale de Lausanne (EPFL) in 2014, 2016, and 2020, respectively. She is currently a postdoctoral researcher at the High Performance Humanoid Technologies Lab (H<sup>2</sup>T), Karlsruhe Institute of Technology, Germany. Her research combines robot learning, optimization, and control with Riemannian geometry. Personal webpage: <http://njaquier.ch>.



**Tamim Asfour** is Professor at the Institute for Anthropomatics and Robotics at KIT, where he is head of the High Performance Humanoid Technologies Lab (H<sup>2</sup>T). His research interest is 24/7 humanoid robotics. Specifically, he studies the mechano-informatics of humanoids as the synergetic integration of informatics, artificial intelligence, and mechatronics into complete humanoid robot systems, which are able to predict and act to perform versatile tasks in the real world. He is the developer and the leader of the development team of the

ARMAR humanoid robot family.