Motion Planning for Human-Robot Collaboration based on Reinforcement Learning

Tian Yu, Student Member, IEEE, and Qing Chang, Member, IEEE

Abstract—Robots are good at performing repetitive tasks in modern manufacturing industries. However, robot motions are mostly planned and preprogramed with a notable lack of adaptivity to task changes. Even for slightly changed tasks, the whole system must be reprogrammed by robotics experts. Therefore, it is highly desirable to have a flexible motion planning method, with which robots can adapt to certain task changes in unstructured environments, such as production systems or warehouses, with little or no intervention needed from non-expert personnel. In this paper, we propose a userguided motion planning algorithm in combination with reinforcement learning (RL) method to enable robots to automatically generate their motion plans for new tasks by learning from a few common tasks saved as primitive actions by kinesthetic human demonstrations. Features of these primitive actions are captured through screw transformation of the endeffector during the task. A mapping method is developed to convert features of primitive actions to new task segments and further used to construct the reward function in RL. A Qlearning algorithm is applied to train the motion planning policy, following which an adaptive motion plan for the new task can be generated or a request for additional primitive actions will be returned if current primitive actions are insufficient for satisfying new task constraints. Multiple experiments conducted on common tasks and scenarios demonstrate that the proposed RL-based motion planning method is effective.

I. INTRODUCTION

In today's manufacturing systems, robot manipulators are usually required to finish repeatable tasks from a starting configuration to a goal configuration while maintaining specific constraints on positions and orientations during the task. However, when task instances change, reprogramming work for robots to satisfy new task constraints is timeconsuming and incapable for non-expert human workers especially when the number of robots in human-robotcollaboration (HRC) is large. With the development of Industry 4.0, there's an increasing demand of robots working adaptively and smartly with humans in recent years [1]. Therefore, how to enable robots to collaborate with non-expert humans and automatically plan motions facing different tasks leave significant challenges to today's smart manufacturing.

To overcome these challenges, we involve both user demonstrations and robot motion planning to generate motion plans for tasks characterized by specific constraints. The use of user demonstration to teach a robot has been studied under the name of Learning from Demonstration (LfD). There is a substantial body of literature on LfD that can be classified based on the method of acquiring the demonstration such as using haptic gloves, image-based sensors, and kinesthetic demonstrations [2]. In this paper, we take advantage of kinesthetic demonstrations since there is no correspondence problem between the kinematic structure of the demonstrating system and the follower robot [3].

There have been a significant number of literatures on LfD that are useful for learning from kinesthetic demonstrations. Some researches use spline functions to decompose demonstrated trajectories to generate new trajectories for manipulation tasks [4]. However, they do not take care of noise in the demonstration which may be significant, especially when the motion information is obtained through human demonstrations. Some nonlinear regression techniques use statistical techniques to incorporate the uncertainty of sensing in the estimation. For instance, a Hidden Markov Model combined with Non-Uniform Rational B-Splines is used in [5] for trajectory approximation. A data-driven approach using Gaussian Mixture Model is adopted in [6]. However, these statistical approaches require multiple similar demonstrations for only one trajectories of the new task. This inefficient usage of user demonstrations would increase the labor of human operators in HRC. The dynamical systems based approach, or dynamical motion primitives (DMPs), on the other hand, can learn from single examples [7]. However, most works assume that there is a dynamical system modeling each degree-offreedom (DoF) of the end-effector, which is not clear when both the position and orientation of the end-effector are relevant to the task. In addition, methods mentioned above are data-driven. The underlying kinetic transformations between each configuration in the user demonstration are not well utilized during the process of LfD.

For robot motion planning, existing algorithms can be divided into joint-space based approaches and task-space based approaches [8], [9]. For a robot manipulator, the joint space or configuration space is the set of all angles the joints can reach. The task space is the set of all end effector configurations of the robot manipulator. Joint space-based motion planning approaches have the advantage that they are (probabilistically or resolution) complete. Thus, in principle, they are guaranteed to give a feasible solution given enough time. However, in practice, they may not return a solution in reasonable time and a non-robotics expert will not be able to handle the situation. On the other hand, task-space based approaches handle task constraints more naturally, but they do

T. Yu is with Department of Mechanical and Aerospace Engineering, University of Virginia, Charlottesville, VA, 22903 USA (ty2yy@virginia.edu).

Q. Chang* is with Department of Mechanical and Aerospace Engineering and Department of Engineering Systems and Environment, University of Virginia, Charlottesville, VA, 22903 USA (<u>ac9nq@virginia.edu</u>)

^{*} Corresponding author

not have any completeness guarantees in the presence of obstacles. Recently, Chakraborty et al. [10] develop a userguided motion planning method in the task space that learns from only one human demonstration to generate motion plans for semantically similar task instances by using the imitated the human demonstration as a guidance given the new goal position. However, task constraints before the current pose blending into the imitated human demonstration are not guaranteed. In addition, the only explicit task constraints like positions and orientations of some other configurations during the task are not considered, which makes it brittle to the change of task.

In this paper, to enable the robot to work with non-expert humans in certain HRC environment and reduce reprogramming when tasks change, we develop an RL-based motion planning method that can enable the robot to learn from one or multiple human demonstrated primitive actions to generate adaptive motion plans for new tasks in such HRC environment. First, primitive actions are recorded in the joint space through kinesthetic demonstrations based on some common tasks such as stacking a block, twisting a screw driver and filling a cup of water. Next, features of these primitive actions are captured by calculating screw transformations throughout demonstrations. Given new task constraints, we formulate the motion planning problem in the task space into a Markov Decision Process (MDP) framework and use the Qlearning method to train a motion planning policy to generate adaptive motion plans in the task space for the new task. Finally, inverse kinematics (IK) is used to calculate corresponding motion plans in the joint space to control the robot to execute learned motion plans.

The main contributions of this paper can be listed as: (1) Propose a method to abstract features of primitive actions; (2) Develop a method that can map features of primitive actions to the new task; (3) Formulate the robot manipulator learning from demonstration and motion planning problem as a reinforcement learning problem in the Markov Decision Process framework; (4) Propose a reasonable reward function based on the evaluation of the deviation between the generated motion plan and the task constraints.

The reminder of this paper is organized as follows: The mathematical preliminaries are introduced in Section II. The learning from demonstration and motion planning problem is stated in Section III. In Section IV, the learning from demonstration and motion planning problem is formulated as an MDP and solved by the Q-learning algorithm. Case studies and conclusions are provided in Section V and VI respectively.

II. MATHEMATICAL BACKGROUND

In this paper, the joint space or configuration space is represented by \mathcal{J} , which is the set of all joint angles of the robot manipulator. *SE*(3) denotes the Special Euclidean group of 3, which represents the task space contains all rigid body motions (i.e., rotations and translations) [11].

To describe the rigid body rotation and translation in SE(3), a dual quaternion representation is adopted in this paper [12]. A dual quaternion **D** is defined as:

$$\boldsymbol{D} = \boldsymbol{d}_r + \frac{1}{2} \epsilon \boldsymbol{d}_t \otimes \boldsymbol{d}_r \tag{1}$$

where $\epsilon \neq 0$, but $\epsilon^2 = 0$. In this definition, d_t and d_r are two quaternions and \otimes is the product of two quaternions. The translation of the rigid body is represented by the quaternion d_t which is denoted as:

$$\boldsymbol{d}_t = (0, \hat{\boldsymbol{t}}) \tag{2}$$

where $\hat{t} = t_x \hat{i} + t_y \hat{j} + t_z \hat{k}$ is the translation vector in *SE*(3). The unit quaternion d_r representing the pure rotation of the rigid body can also be expressed as:

$$\boldsymbol{d}_r = \cos\left(\frac{\phi}{2}\right) + \hat{\boldsymbol{n}}\sin\left(\frac{\phi}{2}\right) \tag{3}$$

where $\hat{n} = n_x \hat{i} + n_y \hat{j} + n_z \hat{k}$ is a unit vector in *SE*(3) representing the rotation axis, and ϕ is the rotation angle. Using this d_r , any vector \hat{v} can be rotated an angle ϕ about the axis \hat{n} by using the quaternion sandwich $d_r \hat{v} d_r^*$ [36], and d_r^* is the conjugate of d_r . For more quaternion manipulations, we refer readers to [12].

III. PROBLEM FORMULATION

In this paper, we consider a human-robot collaboration (HRC) environment, where some common tasks are demonstrated by human workers and recorded as primitive actions. Given constraints of new tasks and environment conditions, e.g., critical configurations of the end-effector and locations of an obstacle, a robot is required to finish these tasks by learning from primitive actions in such HRC environment. We assume that the robot has basic capability to move its end-effector from one configuration to another in the absence of any constraints. Our goal is to develop a method to enable the robot to generate adaptive motion plans for different new tasks to satisfy explicit and implicit task constraints.

A. Specify Manipulation Tasks for the Robot

In this paper, a manipulation task TK is defined as a sequence of n critical configurations based on the task constraints and environment conditions:

$$TK = \{con_1, con_2, \dots, con_n\}$$
(4)

where con_i , i = 1, ..., n, is a tuple of two $\langle \mathbf{P}_i, \boldsymbol{\theta}_i \rangle$. In this tuple, $\mathbf{P}_i = [x_i, y_i, z_i]^T$ is a vector in SE(3) that specifies the position of con_i , $\boldsymbol{\theta}_i = [\theta_{i1}, \theta_{i2}, \theta_{i3}]^T$, is a vector in \mathcal{J} that defines Euler angles of con_i .

B. Features of Primitive Actions

Some common tasks in HRC environment, e.g., fastening screws can be regarded as primitive actions for a certain working scenario. A human can provide a kinesthetic demonstration of a primitive action by holding the arm of the robot and demonstrates a trajectory for accomplishing a given task instance. Throughout this demonstration, the explicit and implicit task constraints embedded in this human demonstrated primitive action can be recorded by many commercial robots such as UR robots in the joint space \mathcal{J} as a time sequence of joint angles $\boldsymbol{\varphi}_{rec}$ as:

$$\boldsymbol{\varphi}_{rec} = \{ \boldsymbol{\varphi}(1), \boldsymbol{\varphi}(2), \dots, \boldsymbol{\varphi}(m) \}$$
(5)

where each $\boldsymbol{\varphi}(t_{rec}) = [\varphi_1(t_{rec}), \varphi_2(t_{rec}), \dots, \varphi_r(t_{rec})]^T, t_{rec} = 1, 2, \dots, m$ represents joint angles of the manipulator, r is the DOF of the manipulator, and t_{rec} is the time sequence of the configurations reached during the motion. Using the forward kinematics mapping $\mathcal{FK}: \mathcal{J} \to SE(3)$, the corresponding human demonstrated primitive action in SE(3) will be obtained as $DP = \{D_1, D_2, \dots, D_m\}$, in which each D_j is a dual quaternion denoted as $D_j = d_{rj} + \frac{1}{2} \epsilon d_{tj} \otimes d_{rj}, j = 1, 2, \dots, m$.

Starting from D_1 , we compute the relative motion with respect to the final end effector configuration for each configuration D_j . Using the dual quaternion representation, the transformation δ_j between the final configuration D_m and every other configuration D_j is:

$$\delta_j = \boldsymbol{D}_j^* \otimes \boldsymbol{D}_m \quad , i = 1, 2, \dots, m-1 \tag{6}$$

where D_j^* is the conjugate dual quaternion of D_j and \otimes is the product of two dual quaternions. This δ_j is also a dual quaternion and can be written as: $\delta_j = d_r^{\delta_j} + \frac{1}{2} \epsilon d_t^{\delta_j} \otimes d_r^{\delta_j}$. Since this transformation δ_j represents the relative translation and rotation of the end-effector between two configurations. It does not change with respect to the current configuration of the end-effector. In this way, all implicit task constraints are embedded in the sequence of δ_j during the motion, which is referred to as the feature of the human demonstration. As such, the feature of the k^{th} primitive action in the *SE*(3) can then be described as:

$$\boldsymbol{H}\boldsymbol{D}_{k}^{\delta} = \left\{ \delta_{1}^{HD_{k}}, \delta_{2}^{HD_{k}}, \dots, \delta_{m-1}^{HD_{k}} \right\}$$
(7)

Suppose that there are h primitive actions in certain HRC environment, the set stores features of all h primitive actions can be denoted as:

$$\boldsymbol{LB} = \left\{ \boldsymbol{HD}_{1}^{\delta}, \boldsymbol{HD}_{2}^{\delta}, \dots, \boldsymbol{HD}_{h}^{\delta} \right\}$$
(8)

Based on the definition of the manipulation task and features of primitive actions, the problem studied in this paper can be described as follows: Given a set of primitive action's features $LB = \{HD_1^{\delta}, HD_2^{\delta}, ..., HD_h^{\delta}\}$ and a new task $TK = \{con_1, con_2, ..., con_n\}$ in the task space SE(3), develop a method to find motion plan **MP** in the joint space \mathcal{J} for the new task by learning from human demonstrations in LB such that the explicit constraints specified in each con_i in **TK** is satisfied. If no **MP** can be found, then a request for additional human demonstrations is made to satisfy all task-relevant constraints in **TK**.

IV. OBTAINING MOTION PLANS TROUGH REINFORCEMENT LEARNING

In order to obtain the motion plan MP in \mathcal{J} given a new task, we want to find out the motion plan MP in SE(3) first since the new task and the set of human demonstrated features have been already established in SE(3). After the motion plan MP in SE(3) is generated, IK can be applied to convert the motion plan in SE(3) to the motion plan in \mathcal{J} .

A. Map human demonstrated features to the new task

Since the new task TK is specified only on some critical configurations con_i , we want to develop a method that learns from human demonstrated primitive tasks to generate an intermediate motion plan between two critical configurations of the new task. The final motion plan MP in SE(3) can consequently be obtained as a sequence of such intermediate motion plans.

First, let $tk_s \subseteq TK$ be a subset or the entire new task TK that can be defined as:

$$\boldsymbol{tk}_{s} = \{con_{i}, con_{i+1}, \dots, con_{v}\}, \ i \ge 1, i \le v \le n \quad (9)$$

Using the dual quaternion representation, the i^{th} configuration of the task tk_s can be written as:

$$\boldsymbol{D}_{i}^{tk_{s}} = \boldsymbol{d}_{ri}^{tk_{s}} + \frac{1}{2} \epsilon \boldsymbol{d}_{ti}^{tk_{s}} \otimes \boldsymbol{d}_{ri}^{tk_{s}}, i = 1, 2, ..., n$$
(10)

The transformation, $\delta_i^{tk_s}$, between the last critical configuration con_n and any other critical configuration con_i is defined as:

$$\delta_i^{tk_s} = \boldsymbol{D}_i^{tk_s^*} \otimes \boldsymbol{D}_n^{tk_s}, i = 1, \dots, n-1$$
(11)

which can also be written as: $\delta_i^{tk_s} = \boldsymbol{d}_r^{\delta_i^{tk_s}} + \frac{1}{2} \epsilon \boldsymbol{d}_t^{\delta_i^{tk_s}} \otimes \boldsymbol{d}_r^{\delta_i^{tk_s}}.$

Therefore, features of the task $\boldsymbol{t}\boldsymbol{k}_s$ can be represented as:

$$\boldsymbol{t}\boldsymbol{k}_{s}^{\delta} = \left\{\delta_{1}^{tk_{s}}, \delta_{2}^{tk_{s}}, \dots, \delta_{n-1}^{tk_{s}}\right\}$$
(12)

Next, suppose that to finish this task segment tk_s , the robot can learn from the feature HD_k^{δ} . We can define a mapping method $mp_k^{\delta}: HD_k^{\delta} \to tk_s^{\delta}$, such that the relative transformation saved in HD_k^{δ} can be scaled and fitted to the task constraints given in tk_s^{δ} . To obtain this mp_k^{δ} , the translation vector $d_t^{\delta j}$ of each δ_j in HD_k^{δ} needs to be aligned to the vector $d_t^{\delta i_k}$ of each $\delta_i^{tk_s}$ in tk_s^{δ} . The rotation from $d_t^{\delta j}$ to $d_t^{\delta i_k^{tk_s}}$ can be represented as a unit quaternion:

$$\boldsymbol{d}_{HD_{k}}^{tk_{s}} = \cos\left(\frac{\phi_{HD_{k}}^{tk_{s}}}{2}\right) + \widehat{\boldsymbol{n}}_{HD_{k}}^{tk_{s}}\sin\left(\frac{\phi_{HD_{k}}^{tk_{s}}}{2}\right)$$
(13)

where
$$\widehat{\boldsymbol{n}}_{HD_k}^{tk_s} = \frac{d_t^{\delta_j} \times d_t^{\delta_i^{tk_s}}}{\left| d_t^{\delta_j} \times d_t^{\delta_i^{tk_s}} \right|}$$
 and $\phi_{HD_k}^{tk_s} = \cos^{-1} \frac{d_t^{\delta_j} d_t^{\delta_i^{tk_s}}}{\left| d_t^{\delta_j} \right| \left| d_t^{\delta_i^{tk_s}} \right|}$.

(

Therefore, each $\boldsymbol{d}_{t}^{\delta j}$ of δ_{j} in $H\boldsymbol{D}_{k}^{\delta}$, j = 1, 2, ..., m-1 can be rotated using the quaternion sandwich $\boldsymbol{d}_{HD_{k}}^{tk_{s}} \boldsymbol{d}_{t}^{\delta j} \boldsymbol{d}_{HD_{k}}^{tk_{s}*}$. After scaling the vector $\boldsymbol{d}_{HD_{k}}^{tk_{s}} \boldsymbol{d}_{t}^{\delta j} \boldsymbol{d}_{HD_{k}}^{tk_{s}*}$ with $\frac{|\boldsymbol{d}_{t}^{\delta j}|}{|\boldsymbol{d}_{t}^{\delta^{1}|}|}$. The final

translation quaternion $d_{tj}^{mp_k^s}$ that maps the relative translation $d_t^{\delta_i}$ in HD_k to tk_s can be obtained as:

$$\boldsymbol{d}_{tj}^{\boldsymbol{m}\boldsymbol{p}_{k}^{s}} = \left(0, \frac{|\boldsymbol{d}_{t}^{\delta_{1}}|}{|\boldsymbol{d}_{t}^{\delta_{1}^{t}}|} \boldsymbol{d}_{H\boldsymbol{D}_{k}}^{t\boldsymbol{k}_{s}} \boldsymbol{d}_{t}^{\delta_{j}} \boldsymbol{d}_{H\boldsymbol{D}_{k}}^{t\boldsymbol{k}_{s}*}\right), j = 1, 2, \dots, m - 1(14)$$

Since the detailed transformation between con_i and con_{i+1} are not specified, we can just use all the intermediate transformation in δ_j for the transformation between con_i and con_{i+1} . Let $d_{rj}^{mp_k^s} = d_r^{\delta_j}$, j = 1, 2, ..., m - 1, the mapping mp_k^s can finally be derived as:

$$\boldsymbol{m}\boldsymbol{p}_{k}^{s} = \boldsymbol{d}_{r}^{\boldsymbol{m}\boldsymbol{p}_{k}^{s}} + \frac{1}{2}\boldsymbol{d}_{t}^{\boldsymbol{m}\boldsymbol{p}_{k}^{s}} \otimes \boldsymbol{d}_{r}^{\boldsymbol{m}\boldsymbol{p}_{k}^{s}}$$
(15)

Let $TK = \{tk_1, tk_2, ..., tk_p\}$, the final motion plan MP in task space for TK can be determined by $MP = \{mp_{k_1}^1, mp_{k_2}^2, ..., mp_{k_p}^p\}$, in which for each $q \in \{1, 2, ..., p\}$, $k_q \in \{1, 2, ..., h\}$. In the following subsections, we will discuss how to formulate the motion problem into a Markov Decision Process (MDP) and obtain this MP in SE(3) using an RL method.

B. MDP Formulation of the Problem

To find out the appropriate motion plan MP that satisfy all task constraints in TK, one has to go through all subsets of TK and evaluate the corresponding motion plans generated by mapping each human demonstrated features in LB to TK, which is a NP-hard problem [13]. Assume that there are x human demonstrations stored and y subsets of the new task TK, the computational complexity for evaluation of the motion plans by searching exhaustively is $O(x^y)$, which would be huge if the number of human demonstrations and the constraints of the new task is large. To solve this NP-hard problem, we formulate the problem into a MDP paradigm which is most common framework for RL that models the sequential decision making in uncertain environments [14].

Before we can apply RL algorithms to obtaining the ultimate motion planning policy π^* , we need to first properly define the three key components of MDP, which ares_t , a_t and r_t at time step t.

Given $LB = \{HD_1^{\delta}, HD_2^{\delta}, ..., HD_h^{\delta}\}$ and $TK = \{con_1, con_2, ..., con_n\}$, the state s_t is defined as:

$$s_t = [CF_t, tk_t] \tag{16}$$

where CF_t is the current configuration of the end-effector at t, $tk_t = \{con_j, con_{j+1}, ..., con_n\}, j \in \{1, 2, ..., n\}$ is the subset of TK and $tk'_t = \{con_j, con_{j+1}, ..., con_k\}, k \le n$ representing the current task the robot is going to satisfy.

The action a_t can be defined as:

$$a_t = [a_1(t), a_2(t), \dots, a_h(t)]$$
(17)

where each $a_i(t)$ is defined as:

$$a_i(t) = \begin{cases} k, & \text{if } HD_i^{\delta} \text{ is mapped to } tk'_t \\ 0, & \text{otherwise} \end{cases}$$
(18)

which is to identify whether the human demonstrated feature HD_i^{δ} is selected to be mapped to the current task tk'_t .

To evaluate the action at t, the reward function r_t is defined as:

$$\begin{aligned} r_t &= \\ \begin{cases} -f\left(\delta_o, \delta_l^{tk'_t}\right), & \forall \delta_l^{tk'_t} \in tk'^{\delta}_t, \exists \delta_o \in HD_i^{\delta} \to \alpha\left(\delta_o, \delta_l^{tk'_t}\right) \leq \Delta_\alpha \\ -\infty, & \text{otherwise} \end{aligned}$$
(19)

where $f\left(\delta_{o}, \delta_{l}^{tk'_{t}}\right) = \sum_{l=j}^{k} \min_{o} \alpha\left(\delta_{o}, \delta_{l}^{tk'_{t}}\right)$ and $\alpha\left(\delta_{o}, \delta_{l}^{tk'_{t}}\right)$ is defined as:

$$\alpha\left(\delta_{o},\delta_{l}^{tk_{t}'}\right) = \min\left\{\left\|\boldsymbol{d}_{r}^{\delta_{o}} - \boldsymbol{d}_{r}^{\delta_{l}^{tk_{t}'}}\right\|, \left\|\boldsymbol{d}_{r}^{\delta_{o}} + \boldsymbol{d}_{r}^{\delta_{l}^{tk_{t}'}}\right\|\right\} (20)$$

which is the Euclidean distance [15] between each δ_o in the feature HD_i^{δ} and each feature $\delta_l^{tK'_t}$ of the current task. Δ_{α} is a tolerance that is set to be 0.5 in this paper.

C. Applying *Q*-learning to Obtain the Optimal Motion *Planning Policy*

The basic idea of Q-learning is that we can define a function Q [14]:

$$Q(s,a) = r(s,a) + \gamma \sum_{s' \in S} p(s'|s,a) v(s',\pi)$$
(21)

such that $v(s, \pi^*) = \max_{a} Q^*(s, a)$. If we know $Q^*(s, a)$, then the optimal policy π^* can be found by simply identifying the action that maximizes $Q^*(s, a)$ under the state s. Starting with arbitrary initial values of Q(s, a), the updating procedure of Q-learning is:

$$Q_{t+1}(s_t, a_t) = (1 - \alpha_t)Q_t(s_t, a_t) + \alpha_t \left[r_t + \gamma \max_a Q_t(s_{t+1}, a_t) \right]$$
(22)

where $\alpha_t \in [0,1)$ is the learning rate and $\gamma \in (0,1)$ is the discount factor. The training process is shown in Algorithm 1. After the training, the ultimate policy π^* is determined as:

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a' \in A(s)} \{Q(s,a')\} \\ 0, & \text{otherwise} \end{cases}$$
(23)

where A(s) is the set of all legal actions at state s in the Q-table Q(s, a). The final motion plan **MP** is generated following the Algorithm 2.

Algorithm 1 Training of the RL-based Motion Planner

Input: *TK*, *LB*, ϵ , γ Output: Q(s, a)Initialize Q(s, a) randomly Initialize t = 0Initialize s_0 with $CF_0 = D_1^{tk_1}$ and $tk_0 = TK$ For episode = 0,1, ...,100 do While the last con_n of *TK* is not reached do Choose $a_t = \arg \max_{a \in A(s_t)} Q(s_t, a)$ using policy derived from Q(s, a) (e.g. ϵ -greedy) Take action a_t , observe r_t Update $Q_{t+1}(s_t, a_t)$ based on Eqn. (22) End While End For

Output Q(s, a)

Algorithm 2 Generating Motion Plan in SE(3)

Input:
$$TK$$
, LB , ϵ , γ , $Q(s, a)$

Output: MP

Initialize s_0 with $CF_0 = D_1^{tk_1}$ and $tk_0 = TK$

For t = 0, 1, ..., T do

Find legal action list $A(s_t)$ from $Q(s_t, a) \rightarrow a \in A(s_t)$

Find the optimal action as $a_t = \arg \max_{a \in A(s_t)} Q(s_t, a)$

Map HD_i to tk'_t according to a_t

 $MP_t \leftarrow \boldsymbol{D}_n^{tk'_t} \otimes \boldsymbol{mp}_i^{t*}$

End For

 $\boldsymbol{MP} \leftarrow \{MP_0, MP_1, \dots, MP_T\}$

Output MP

V. CASE STUDY

In order to validate effectiveness of the proposed method in generating motion plans for new tasks, multiple experiments are conducted on the UR5e platform. In this case study, two performance metrices are considered: (1) The accumulated reward of the motion plan in SE(3); (2) The successful rate of applying the motion plan in \mathcal{J} for new tasks. From the case study, two significant results can be concluded: (1) The proposed method is effective in combining different features of human demonstrations to generate motion plans for the new task; (2) The proposed method is effective in requesting additional human demonstrations if no features of primitive actions are semantically similar to the new task.

A. Record Primitive actions in J

As shown in Fig. 1, three primitive actions are recorded in \mathcal{J} through kinesthetic demonstrations, including one screwing task, one filling task and one stacking task.



Figure 1. Kinesthetic demonstrations of 3 most common tasks. (a) Screwing task. (b) Filling task. (c) Stacking task.

B. Training the RL-based Motion Planner in SE(3)

In order to obtain a general motion planning policy in SE(3) for the HRC scenario, we implement Algorithm 1 to train a Q-table initialized with random Q values in Matlab using a 4-core 4.0GHz Intel Core i7 processor. 20 new tasks with four critical configurations of each that cover all features are used during the training. Positions of these critical configurations are randomly generated within a $50 \times 50 \times$

50 cm³ workspace. Corresponding Euler angles of each critical configuration are also randomly selected from a set $\{-\pi, -\pi/2, 0, \pi/2, \pi\}$. The total training episode for each new task is set to be 100. The total computation time is 1927.42 seconds.

To monitor the training process, accumulated rewards for each new task are recorded every two iterations. The average accumulated reward for all 20 new tasks is shown in Fig. 2. Although the training rewards are noisy before 50 episodes, the underlying trend is that the rewards are increasing with training episodes. It can be observed that the reward reaches a steady level after around 50 episodes. This indicates that a steady motion planning policy in SE(3) that can map appropriate features of human demonstrations to new tasks with semantically similar features is generated for the assemble and loading/unloading scenario.



Figure 2. Training process for motion plans in SE(3)

C. Evaluation of the Trained Motion Planning Policy

To evaluate the performance of the trained motion plan policy in SE(3) for the HRC scenario, two new tasks, namely, a filling-and-pouring task and a passing task, are used as examples to demonstrate the method. The trained Q-table from is used as the input to Algorithm 2 to generate motion plans for new unseen tasks in SE(3). IK is used to calculate the final motion plan in \mathcal{I} . For each task, 20 experiment trials are conducted. For each trial, if the final motion plan in \mathcal{I} is collision-free and satisfy all task constraints, such trial would be regarded as a successful one. The successful rate for each new task is recorded to evaluate the motion planning policy.

Filling-and-Pouring Task: In this task, the end-effector is required to fill water to the blue cup (shown in Fig. 3), then reach a goal position above the white cup (shown in Fig. 3), and finally pour water to the white cup. The location of the white cup is on the surface of a desk within a workspace of $20 \times 20 \ cm^2$. We can specify 4 critical configurations based on the described task. Sample critical configurations from con_1 to con_4 are presented in Table 8, where $-\pi \leq \gamma \leq \pi$, $-0.5 \le x \le 0.7$, $-0.2 \le y \le 0$ for different trials. For 20 experiment trials, a successful rate of 80% is observed. We use the final motion plan in \mathcal{I} for one trial as an example as shown in Fig. 3 to illustrate the result. It is noticed that the features of 3 human demonstrated tasks, namely filling, stacking, and screwing, are learned and mapped to the segment between con_1 and con_2 , the segment between con_2 and con_3 , and the segment between con_3 and con_4 , respectively. In this experiment, the proposed method can identify and compose the appropriate features of primitive actions to perform a new task.

TABLE I. CRITICAL CONFIGURATIONS OF THE FILLING-AND-POURING TASK

	con_1	con_2	con_3	con_4
P	(-0.4,-0.1,0)	(-0.5,0.1,0.1)	(<i>x</i> , <i>y</i> , 0.1)	(<i>x</i> , <i>y</i> , 0)
θ	(0,- <i>π</i> ,0)	$(0, -\pi/2, \gamma)$	$(0, -\pi/2, \gamma)$	$(\pi/2,0,-\pi/2)$
	A STER	con a		on2

Figure 3. Final motion plan for the filling-and-pouring task.

Passing Task: In this task, the end-effector is required to pass the screw driver to the human in a specific orientation. By applying the same trained general motion plan policy, a zero successful rate is observed in the experiment, which indicates additional human demonstrations are needed. A closer examination reveals that none of the three features of primitive actions is semantically similar to the feature of the task segment between con_2 and con_3 , which requires a 90-degree rotation about its body-fixed x-axis clockwise. Therefore, additional human demonstration is requested for this feature.

TABLE II. CRITICAL CONFIGURATIONS OF THE PASSING TASK

	con_1	con ₂	con_3
Р	(-0.5,-0.1,0.2)	(-0.5,0.1,0.4)	(-0.7, 0.2, 0.4)
θ	$(\pi, 0, \pi)$	$(\pi/2, 0, -\pi)$	$(-\pi/2, 0, -\pi/2)$

With this additional human demonstrated primitive action shown in Fig. 4 (a) added, the motion planning policy is retrained using Algorithm 1 with the Q-table trained in section 6.2. As shown in Fig. 4 (b), the accumulated reward reaches a steady value after around only 8 iterations. Then by applying the newly trained policy, the motion plan is generated as shown in Fig. 4 (c). All 20 trials are witnessed successful. The result shows that the features of the human demonstrated screwing task is learned and mapped to the task segment between con_1 and con_2 . The feature of the newly added human demonstration is learned and mapped to the task segment between con_2 and con_3 .

VI. CONCLUSION AND FUTURE WORK

In this paper, we present a method for incorporating human demonstrations in RL-based motion planning. By learning from human demonstrated features of primitive actions, the task-space RL-based motion planner can effectively generate motion plans for new tasks or request additional human demonstrations for tasks if the features of all primitive actions are insufficient to finish the task. Using inverse kinematics (IK), motion plans in the joint space can be obtained and successful trials are able to be achieved in real-world experiments. In future work, we plan to extend our proposed method to mobile manipulators. We also plan to integrate with the human motion perception, recognition, and prediction for realistic implementation.



Figure 4. Outline for the passing task. (a) Additional human demonstration. (b) Training process with newly added human demonstration. (c) Final motion plan for the pass task.

ACKNOWLEDGMENT

This work was supported by the U.S. National Science Foundation (NSF) Grant. CMMI1853454.

REFERENCES

- S. El Zaatari, M. Marei, W. Li, and Z. Usman, "Cobot programming for collaborative industrial tasks: An overview," *Rob. Auton. Syst.*, vol. 116, pp. 162–180, 2019.
- [2] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Rob. Auton. Syst.*, vol. 57, no. 5, pp. 469–483, 2009.
- [3] Z. Zhu and H. Hu, "Robot learning from demonstration in robotic assembly: A survey," *Robotics*, vol. 7, no. 2, p. 17, 2018.
- [4] J.-H. Hwang, R. C. Arkin, and D.-S. Kwon, "Mobile robots at your fingertip: Bezier curve on-line trajectory generation for supervisory control," in *Proceedings 2003 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2003)(Cat. No. 03CH37453)*, 2003, vol. 2, pp. 1444–1449.
- [5] J. Aleotti and S. Caselli, "Robust trajectory learning and approximation for robot programming by demonstration," *Rob. Auton. Syst.*, vol. 54, no. 5, pp. 409–413, 2006.
- [6] S. Calinon, F. Guenter, and A. Billard, "On learning, representing, and generalizing a task in a humanoid robot," *IEEE Trans. Syst. Man, Cybern. Part B*, vol. 37, no. 2, pp. 286–298, 2007.
- [7] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal, "Dynamical movement primitives: learning attractor models for motor behaviors," *Neural Comput.*, vol. 25, no. 2, pp. 328–373, 2013.
- [8] L. E. Kavraki, P. Svestka, J.-C. Latombe, and M. H. Overmars, "Probabilistic roadmaps for path planning in high-dimensional configuration spaces," *IEEE Trans. Robot. Autom.*, vol. 12, no. 4, pp. 566–580, 1996.
- [9] O. Khatib, "A unified approach for motion and force control of robot manipulators: The operational space formulation," *IEEE J. Robot. Autom.*, vol. 3, no. 1, pp. 43–53, 1987.
- [10] R. Laha, A. Rao, L. F. C. Figueredo, Q. Chang, S. Haddadin, and N. Chakraborty, "Point-to-point path planning based on user guidance and screw linear interpolation," in *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*, 2021, vol. 85451, p. V08BT08A010.
- [11] J. M. Selig, Geometric fundamentals of robotics, vol. 128. Springer, 2005.
- [12] L. F. da C. Figueredo, "Kinematic control based on dual quaternion algebra and its application to robot manipulators," 2016.
- [13] D. P. Bovet, P. Crescenzi, and D. Bovet, *Introduction to the Theory of Complexity*, vol. 7. Prentice Hall London, 1994.
- [14] D. Silver, "Reinforcement Learning and Simulation-Based Search," 2009.
- [15] L. Kavan, S. Collins, C. O'Sullivan, and J. Zara, "Dual quaternions for rigid transformation blending," *Trinity Coll. Dublin, Tech. Rep. TCD-CS-2006-46*, 2006.