

# Towards Learning to Handle Deviations using User Preferences in a Human Robot Collaboration Scenario

Sharath Chandra Akkaladevi<sup>1,2,\*</sup>, Matthias Plasch<sup>1</sup>, Christian Eitzinger<sup>1</sup>, Srinivas Chowdhary Maddukuri<sup>1</sup>, and Bernhard Rinner<sup>2</sup>

<sup>1</sup>Profactor GmbH, Im Stadtgut A2, Steyr-Gleink, 4407 Austria

<sup>2</sup>Institute of Networked and Embedded Systems, Alpen-Adria-Universität Klagenfurt, Klagenfurt, Austria

\*Corresponding Author E-mail: sharath.akkaladevi@profactor.at

**Abstract.** In a human robot collaboration scenario, where robot and human coordinate and cooperate to achieve a common task, the system could encounter with deviations. We propose an approach based on Interactive Reinforcement Learning that learns to handle deviations with the help of user interaction. The interactions with the user can be used to form the preferences of the user and help the robotic system to handle the deviations accordingly. Each user might have a different solution for the same deviation in the assembly process. The approach exploits the problem solving skills of each user and learns different solutions for deviations that could occur in an assembly process. The experimental evaluations show the ability of the robotic system to handle deviations in an assembly process, while taking different user preferences into consideration. In this way, the robotic system could both benefit from interaction with users by learning to handle deviations and operate in a fashion that is preferred by the user.

**Keywords:** Learning to handle deviations, user preferences, human robot collaboration, ontology based knowledge representation

## 1 Introduction

Human robot collaboration (HRC) allows to combine the cognitive strength of humans together with the physical strength of robots and can lead to numerous applications [6]. For a close collaboration with humans, the robotic systems should attribute meaning to beliefs, goals and desires, collectively called "mental models" of humans during a particular task [12]. This would not only allow the robotic system to understand the actions and expressions of humans within an intentional or goal-directed architecture [16], but also for the human operators to better understand the capabilities of the robotic system [12]. In context of HRC, user preferences communicate the human's desire to achieve a goal (based on the belief that the goal is possible), either by themselves or in collaboration with the robotic system. User preferences can be communicated implicitly (via the actions done by the human) or explicitly (with the help offline knowledge representation). User preferences aid the robotic system in

creating and understanding the perspective of the human and collectively forms the basis for building the "mental model" of the human.

A detailed survey on human robot interaction, emerging fields and applications is given in [7][4]. In this paper we present an HRC approach capable of handling deviations, where the robot interacts with the human by considering the 'preferences' of the human. By considering user preferences, the robotic system can take actions from the perspective (mental model) of the human to facilitate a `natural' interaction. HRC should deal with the non-deterministic factor of the environment with human at its center, but it could also benefit from the cognitive and problem solving skills of the human. In recent works [8, 18, 21], it was shown that given the presence of the human, the robotic system can take advantage of human feedback and learn a given task in short time using Interactive Reinforcement Learning (IRL) techniques [18]. Inspired by this success, we propose an extension to the Interactive Reinforcement Learning (eIRL), which would allow the robotic system to take the aid of the human (asking which action to take next, receiving feedback) to learn to cope with deviations/novel situations as they occur. The advantage of allowing the user to teach the robot in coping with deviations is that each user has an own way of solving a given problem (communicated via user preferences). This allows the robot to learn different possibilities (when they exist) to deal with a deviation and use this knowledge to better collaborate and cooperate with the human. This also helps the robotic system to better personalize its behavior to a given user and their corresponding preferences.

In order to showcase the deviation handling capability of the proposed approach, an assembly process use case that deals with assembling objects into a box in collaboration with the human is chosen. This use case consists of four real-world objects, namely *heater*, *tray*, *ring* and *base* as shown in Figure 1a. The goal of the use case is to assemble these objects by placing each of them inside a box. An example execution of pick and place operation of the *heater* object is shown in Figure 1b. It is important to note that the manipulation possibility of each object either by the robot or the human is dependent on the construction (e.g., heavy, objects that require complex grasp - *base*) and configuration of the objects in the workspace.

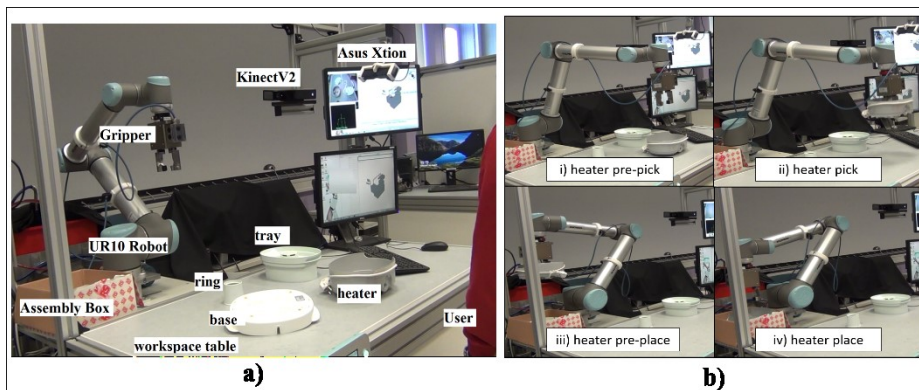


Fig. 1. (a) Use Case setup; (b) Pick and Place operation of the heater object

In this article, an HRC architecture that is capable of dealing with uncertainties and deviations of task executions using eIRL is presented. The architecture exploits the assistance of the human (who interacts with the robot to teach it) in the HRC loop, to learn and recover from deviations/novel situations that can occur during a task execution, using an interactive reinforcement learning algorithm. The approach where the robotic system:

- assesses the feasibility of executing the goal (to a certain level) before performing it
- **proactively suggests** possible course of action/s while considering the robot state (proprioception), the environment and the **preferences of the user**
- **enables technically unskilled users** to teach and customize the robot's behavior to their preferred manner
- **provides an interface** to communicate and learn from the human the necessary/alternative steps to take in case of deviations

The emphasis is on dealing with the deviations that occur in the action planning phase of the task. The low-level deviations that occur while carrying out the motion planning (path planning, collision avoidance) of the robot, while executing the task are not in the focus of this work.

The remaining part of the paper is organized as follows: In Section 2, a brief description of the existing state of the art approaches that deal with RL to learn a complete task are presented. Then in Section 3, the functional components and the architecture of the extended IRL is described. The theoretical background and the algorithmic description of eIRL is given in Section 4. Finally, the experimental setup and evaluation of the proposed eIRL is discussed in Section 5 followed by concluding remarks and possible future work described in Section 6.

## 2 Related Work

In order for the collaborative robots to work hand in hand with human operators the robotic system should be able to deal with complex and continuously changing environments. Thorough offline modeling of the environment and task conditions is not a feasible solution. Hence the robotic system should learn to respond to the environment by performing actions in order to reach a goal. The problem of agents learning in complex real-world environments is dealt with by numerous approaches [4][5][7]. Reinforcement learning (RL) [19] is one such popular approach that deals with learning from interaction, to teach the agent how to behave (which actions to perform when) in order to complete a task. The main aim in RL is to maximize a cumulative reward that is attained by taking some actions in the environment.

In RL, the agent learns over discrete time steps by interacting with its environment and gaining experience about the outcome. However, in order to reach an optimal policy (the set of actions that lead to the maximum reward), the RL approach requires substantial interaction with the environment. Depending on the nature of the task, RL results in a memory intensive storage of all state action pairs [9]. Another disad-

vantage of RL is its slow convergence towards a satisfactory solution. Despite such drawbacks, recent works on learning have shown considerable interest in using RL [8][11][18][21]. The basic idea is to use the human teacher in the loop to provide feedback and hence speedup the convergence time in RL for reaching an optimal policy. In conventional RL, a reward is a positive or negative feedback for being in the current state (or for a particular state-action pair). Thomas et al [21] introduced a run-time human feedback as a reward to the IRL approach and argue that this approach is beneficial for both the human teacher and the learning algorithm. Their study suggests that users (human teachers) employ the feedback as a single communication channel for various communicative intents-feedback, guidance, and motivation. Inspired by this approach, Suay et al [18] study how IRL can be made more efficient for real-world robotic systems. In this approach, the user provides rewards for preceding actions of the robot and additionally provides guidance for subsequent actions. They further show that this guidance reduces the learning time of the robot and that it is more evident when large state-space size (number of interactions in the environment for the robot) is considered.

Knox et al [10] propose a framework to train an agent manually via evaluative reinforcement, using real-valued feedback on its behavior from a human trainer. This allows the human trainer to interactively shape the agent's policy (interactive shaping) and thereby, directly modify the action selection (policy) mechanism of the IRL algorithm. Rather than influencing the policy indirectly through a reward, Griffith et al [8] use the user feedback in making a direct statement about the policy itself. Knox et al [11] extended their previous work in [10] by applying the TAMER framework to real-world robotic system. Rozo et al [14] [15] propose an approach for a human robot interactive task, where the robot learns both the desired path and the required amount of force to apply on an object during the interaction.

In this article, the following contributions are made to the state of the art:

- an extension to the Interactive Reinforcement learning algorithm (eIRL), that enables the robotic system to learn from the human, to deal with deviations.
- at each instant the robotic system encounters a deviation, it proactively suggests a list of actions possible by the robot in its current state. This is achieved through efficient representation of the task knowledge using web ontology web language (OWL), see Section 3 for more details. The advantage Learning to Handle Deviations for Human Robot Collaboration 5 is that even untrained users can effortlessly teach the robot (how to deal with deviations) by choosing an appropriate action from the list provided.
- unlike the other approaches, the human feedback is divided into two: a) feedback to choose an action policy from the list of options available b) like/dislike feedback to function as a reward for an action executed by the robot, when a deviation is encountered.

### 3 eIRL Architecture

In order to collaborate and cooperate with the human in HRC scenario, the following components are developed within the extended IRL (eIRL) architecture as shown Figure 2: a) **perception and interpretation capabilities**: to understand and interpret the current state of the environment from sensor data b) **knowledge representation**: to represent different aspects of the task carried out in the HRC environment that include representing the abilities and activities of the human, interplay between human activities and object configurations, robot's own self with respect to the task, etc. c) **learning and reasoning**: the robotic system is equipped with learning capabilities that allows the accumulation of knowledge over time to enhance in particular its abilities to perceive the environment, to make decisions, to behave intelligently, and to interact naturally with humans. The robotic system also reasons about its current state to plan its own actions accordingly to aid in completing the task d) **action planning**: to generate plans for future actions to achieve the given task (goals). This includes task planning, scheduling and observation planning, as well as planning under uncertainty for an efficient human robot collaboration. These plans generated need to be carried out in real world where the robot plans its path (path/navigation planning) and manipulates the environment accordingly. Note that the learning, reasoning and the action planning components form the heart of the eIRL architecture. e) **functional component layer**: provides the necessary low-level sensory information like real-time 3D tracking of objects (implemented similar to [2]), current action of the human [1] together interfaces to robotic manipulation and a GUI-interface for human robot interaction.

**Knowledge Representation** in the proposed architecture is based on KnowRob [20]. The main reason for choosing KnowRob is that it provides the following knowledge processing features: a) mechanisms and tools for action centric representation b) automated acquisition of grounded concepts through observation and experience c) reasoning about and managing uncertainty, and fast inference. The knowledge is represented using ontology (description logics) in the Web Ontology Language (OWL) and SWI Prolog is used for loading, accessing and querying the ontology. The representation consists of two levels: *Classes* that contain abstract terminological knowledge (type of objects, events and actions - taxonomic fashion) and *Instances* which represent the actual physical objects or the actions that are actually performed. The link between *Classes* and *Instances* is given in *Properties*, which defines if an  $Agent \in \{Human, Robot\}$  can perform a particular action (defined in *Classes*) on/with a  $Target \in \{Objects, Robot, Human\}$ . A Relation denoted by a triple  $\langle Agent, Property, Target \rangle$  defines a particular aspect of the assembly task. For example,  $\langle Robot, pick, ring \rangle$  conveys that the robot (is capable and) should pick the ring object. A detailed description of KnowRob and its features are given in [20].

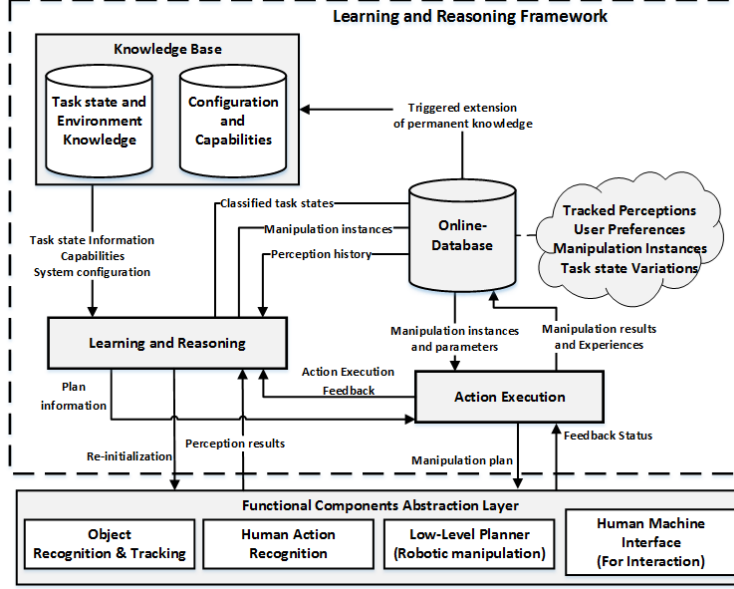


Fig. 2. Architecture diagram of the extended Interactive Reinforcement Learning approach

## 4 eIRL Algorithm Description

Reinforcement Learning (RL) is an area of machine learning that defines a class of algorithms that enable a robot to learn from its experience. Reinforcement learning in our case is defined using the standard notation of the Markov Decision Process (MDP). In a MDP, any state  $s_{t+1}$  occupied by the robot is a function of its previous state  $s_t$  and the action taken at, in other words  $s_{t+1} = f(s_t, a_t)$ . A MDP is denoted by the 5-tuple  $\langle S, A, P, R, \gamma \rangle$ , where the set of possible world states is defined by  $S$ , and  $A$  denotes the set of actions available to the agent in each state. The probability function  $P : S \times A \rightarrow Pr[S]$ , describes the transition probability of State  $s_t$  to State  $s_{t+1}$ , when an action  $a_t$  is performed on State  $s_t$ . The reward function  $R : S \times A \rightarrow R$ , and a discount factor  $\gamma$ , where  $0 \leq \gamma \leq 1$ . Together  $P$  and  $R$  describe the dynamics of the system. The goal of a RL algorithm is to find an approximation function  $Q : S \times A \rightarrow R$ , that defines an optimal policy  $\pi$ . Where,  $Q$  maps the state-action pairs to the expected reward and the optimal policy  $\pi : S \rightarrow A$  maximizes the expected reward. In other words,  $\pi$  specifies the best possible action to perform in a given state in order to gain a maximum reward.

In this paper, Q-learning RL [23] is used as a basis for the learning algorithm. As mentioned earlier, IRL approaches [18][21] showed the success of learning a complete task using interactive RL algorithm, Q-learning. Using this as basis, we extend IRL to be used in handling deviations in an assembly process as they occur during execution as shown in Figure 3. The main idea is to integrate the human operator as advisor into the learning process. In this way the user can teach the system, how to deal with deviations during assembly process execution. It is assumed that the optimal

action selection policy, say  $\pi^*$ , to perform the assembly task, is already known (using [18][21]). In this section we describe the algorithm that modifies Q-learning, to handle deviation using user preferences. In this context, two types of user feedback are considered a) the optimal action policy selection by the user during the deviation - treated as the user preference b) the reward provided by the user for the action taken by the robotic system during the deviation - the like/ dislike option.

During a normal assembly process, the robotic system follows the already learned optimal policy  $\pi^*$  and performs actions  $a_{ii}$  accordingly based on its current state  $s_t$  to progress to the next state  $s_{t+1}$ . As described earlier, the robotic system is enabled with state of the art action recognition and object tracking perception capabilities coupled with knowledge representation framework to comprehend the current status of the assembly task. After completing an action  $a_{ii}$  the robotic system observes the resulting state, if the resulting state is the expected next state  $s_{t+1}$ , the next optimal action policy in  $\pi^*$  is carried out and so on until the end. However, due to the dynamic nature of the environment, though an action is chosen by the robotic system, the task execution might not be successful and is called a deviation. A deviation could occur due to a) the chosen action cannot be performed - e.g., object out of reach, object not available b) expected resulting state  $s_{t+1}$  was not observed after the action  $a_{ii}$  is performed. A typical example could be an object was not graspable due to sensor error or the object configuration. It could also happen that the object might be out of reach of the robot. In such cases, it is not possible for the robotic system to continue with the assembly task as the learned optimal policy action cannot be completed.

Let us say action  $a_{ii}$  is an optimal policy action (given  $\pi^*$ ) taken in state  $s_t$  and a deviation  $\Delta_{ii}$  has occurred. In this case, the robotic system (depending on the current state, its capabilities, objects, knowledge base) proactively suggests a set of deviation handling actions  $a_{ti,j}^\Delta$  possible in the current state  $s_t$ . Since deviations are special cases that occur during the assembly process, the original optimal policy should not be directly affected. Therefore, a deviation Q vector  $\vec{Q}_{ti}^\Delta$  is associated with this deviation  $\Delta_{ii}$  that occurred for this state-action pair  $\{s_t, a_{ii}\}$  and its corresponding entry  $Q_{t,i}^*$  in the Q table. If no vector  $\vec{Q}_{ti}^\Delta$  existed previously for  $\Delta_{ii}$ , then it is created and its values  $Q_{ti,j}^\Delta$  are uniformly initialized to zero. This is generally the case, if this deviation was occurring for the first time. The robotic system then waits for the user to select an action policy ( $a_{i,j} \in A_{ti}^\Delta$ ) suitable for the deviation, from the list of suggested actions. In order to suggest the list of possible deviation handling actions  $A_{ti}^\Delta$ , the robotic system observes the current status in the assembly process, available objects and robot's capabilities and then proposes a list of actions to progress the assembly process. Once the user selects an action policy, the action is carried out and the resulting state  $s_{next}$  is observed. Whenever, a deviation occurs the robotic system learns the optimal action policy in case of that deviation myopically i.e.,  $\gamma = 0$  [11]. In other words, it is assumed that no matter what deviation action policy is chosen by the user, it will progress the assembly process from the deviation to a known state  $s_{known}$ , from which an optimal state-action pair exists. This kind of learning from human feedback eliminates any exploration and results in a deviation handling policy that is only as good as the decision made by the user.

---

```

1: Initialize Q table  $Q^*$  containing the action values  $Q_{t,i}^*$ , where  $1 \leq t \leq |S|$  and
    $1 \leq i \leq |A_t|$  are the indices of states  $s_t$  and actions  $a_{ti}$  possible at the given state.  $Q^*$ 
   leads to an optimal assembly process execution policy  $\pi^*$ , maximizing the reward

2: for each execution episode of the assembly process do
3:   Set assembly process to initial state  $s_t | t = 1$ 
4:   for each execution step do
5:     Execute action  $a_{ti}^*$  in state  $s_t$  following policy  $\pi^*$ 
6:     if no deviation occurred while executing  $a_{ti}^*$  then
7:       Observe resulting state  $s_{t+1}$ 
8:       Q table  $Q^*$  is unchanged
9:     else
10:      Propose deviation handling actions  $a_{ti,j}^\Delta \in A_{ti}^\Delta | j \in \mathbb{N}$ 
11:      if Action value  $Q_{t,i}^* \not\sim \vec{Q}_{ti}^\Delta$  then
12:        Initialize deviation Q vector  $\vec{Q}_{ti}^\Delta$  with  $Q_{ti,j}^\Delta = 0$ 
13:        Associate  $Q_{t,i}^*$  with  $\vec{Q}_{ti}^\Delta$ :  $Q_{t,i}^* \sim \vec{Q}_{ti}^\Delta$ 
14:        if  $\vec{Q}_{ti}^\Delta$  holds the initial values, i.e.  $\max Q_{ti,j}^\Delta = 0$  then
15:          User shall select a deviation handling action  $a_{ti,k}^\Delta | 1 \leq k \leq |A_{ti}^\Delta|$ 
16:        else
17:          Take handling action  $a_{ti,k}^\Delta \leftarrow \arg \max_{a^\Delta \in A_{ti}^\Delta} Q_{ti,j}^\Delta$ 
18:        Observe resulting state  $s_{t+1}$ , if  $a_{ti,k}^\Delta$  was unsuccessful then go to ln.14
19:        Receive user feedback (like/dislike), determining the reward  $r_{t,k}^\Delta$ 
20:        for each  $Q_{ti,j}^\Delta \in \vec{Q}_{ti}^\Delta$  do
21:           $Q_{ti,j}^\Delta \leftarrow Q_{ti,j}^\Delta + \alpha \left[ r_{t,j}^\Delta + \gamma \max_{1 \leq l \leq |A_{t+1}|} Q_{t+1,l} - Q_{ti,j}^\Delta \right] | \alpha = 1, \gamma = 0$ 
22:          if  $j = k \implies a_{ti,j}^\Delta = a_{ti,k}^\Delta$  then
23:            if userFeedback = like then
24:               $r_{t,j}^\Delta \leftarrow r_{max}$ 
25:            else
26:               $r_{t,j}^\Delta \leftarrow 0$ 
27:          else
28:             $r_{t,j}^\Delta \leftarrow 0$ 
29:         $s_t \leftarrow s_{t+1}$ 

```

---

#### Comments

Ln 10:  $j$  is the index of possible deviation handling actions  $a_{ti,j}^\Delta$   
Ln 11: Each action value  $Q_{t,i}^*$  can be associated with a deviation Q vector  $\vec{Q}_{ti}^\Delta$   
Ln 11: The symbol  $\not\sim$  means "not associated with"  
Ln 12:  $Q_{ti,j}^\Delta$  is the action value for applying the deviation handling action  $a_{ti,j}^\Delta$  to  $s_t$ , after encountering a deviation in action  $a_{ti}^*$   
Ln 18: It is assumed that the resulting state  $s_{t+1}$  is a known state of the assembly process  
Ln 23: The reward  $r_{max} \in \mathbb{R}^+$  is given, if the deviation handling result is *liked* by the user

---

**Fig. 3.** Algorithm to Handle deviations based on Q-Learning

Upon successful completion of the suggested action policy ( $s_{next} = s_{known}$ ), the robotic system asks for a *reward feedback*. The user at this stage can either *like* or *dislike* the robot's performance, as shown in Figure 4. If the action is *liked* then the respective Q value for that action in the deviation  $\vec{Q}_{ti}^\Delta$  vector is set to a maximum reward value and the other action entries are set to zero. In case of *dislike*, all the Q values in  $\vec{Q}_{ti}^\Delta$  are set to 0. The execution is then continued with the optimal policy. If the deviation handling action was not successfully executed ( $s_{next} \neq s_{known}$ ), the robotic system goes

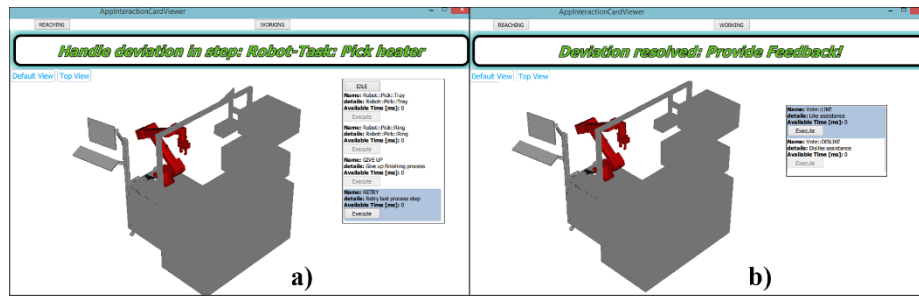


back and suggests possible deviation actions again. In this way the robotic system learns how to deal with the deviation that occurred for that state-action pair.

At a later point in time, if the robotic system encounters a deviation  $\Delta_{ii}$ , there already exists a  $\vec{Q}_{ii}^A$  with initialized values. The robotic system then directly chooses the action policy  $a_{ti,j}^A$  in the  $\vec{Q}_{ii}^A$  that leads to a maximum reward value. After performing the deviation handling action, the robotic system asks the user for reward feedback. If the user then *likes* the action performed by the robot, the Q values of the  $\vec{Q}_{ii}^A$  are updated accordingly and the assembly process proceeds as usual. Since the occurrence of these deviations are separated in time (a deviation can occur today and is repeated after an year), it is possible that user might like to teach the robot a different action policy  $a_{ti,diff}^A$  for handling that deviation  $\Delta_{ii}$ . In such cases, the user *dislikes* the action performed by the robot which allows the robot to learn a new deviation action policy  $a_{ti,diff}^A$  for that deviation  $\Delta_{ii}$ .

## 5 Experimental Setup

As shown in Figure 1a, the experimental setup consists of a UR10 [22] robotic arm with 6 degrees of freedom, a SCHUNK 2 finger electric parallel gripper and two commercially available RGB-D sensors each equipped with object tracking and action recognition functionalities respectively. As mentioned in Section 3, the functional component layer with the help of sensors and actuators as shown in Figure 1a provides the reasoning and learning module with necessary perception (object tracking and action recognition) and actuation (robot and gripper state combined to deduce if an action execution was successful) information.



**Fig. 4.** Example of interaction cards provided while handling deviations in the assembly process (a) shows the list of deviation handling actions possible when robot encountered a deviation while Pick&Place heater action and (b) shows the like/dislike feedback interaction card presented to the user after successful handling of the deviation

The assembly process use case presented in this paper concerns itself with assembling the given objects (*heater, base, tray, ring*) in to a box. Depending on the capabilities of the *Agent*  $\in \{Human, Robot\}$  and the object properties (configuration, construction) defined in the knowledge representation an optimal policy of a normal assembly process is defined as given in Figure 5 and is assumed to be known. Figure 5

shows the optimal policy sequence (Step 1, Step 2, Step 3, Step 4 respectively) that expresses 'what' action is required to be performed by 'which' *Agent* and on 'which' *Target*. In this section, the ability of the system to learn, how to deal with deviations that occur during online execution (which were not encountered during training) of the assembly process is evaluated. During a normal assembly, Step 1 to 4 are executed in a sequence following the optimal policy and thereby gaining a maximum reward. The *Pick&Place* action includes localizing the object, reaching for the object, grasping and lifting it and then placing it in the pre-assigned assembly box as shown in Figure 1b. Note that, as we deal with a human robot collaboration scenario, some steps needs to be carried out by the user, in this case Step 1. The robotic system waits for the user to complete Step 1, observed with the help of embedded object tracking and action recognition functionalities. Once Step 1 is complete the robot proceeds with the next steps.

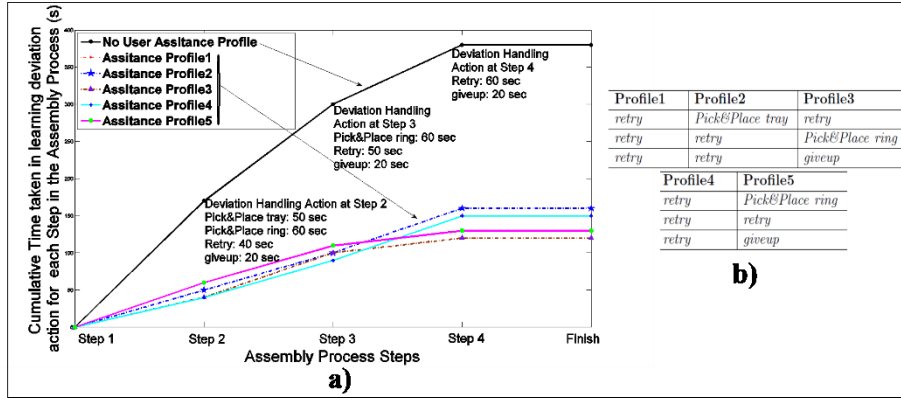
Assembly Process Step < <i>Agent, Property, Target</i> >	Possible Deviation Handling Actions
Step 1: User <i>Pick&amp;Place base</i>	Note: only user can perform this action
Step 2: Robot <i>Pick&amp;Place heater</i>	<i>retry; Pick&amp;Place tray; Pick&amp;Place ring; giveup</i>
Step 3: Robot <i>Pick&amp;Place tray</i>	<i>retry; Pick&amp;Place ring; giveup</i>
Step 4: Robot <i>Pick&amp;Place ring</i>	<i>retry; giveup</i>

**Fig. 5.** Optimal action policy of the assembly process and the respective deviation handling actions available at each step

A deviation can occur at any step in the assembly process and the possible deviation handling actions at each step are given in Figure 5. These deviation handling actions are presented as interaction cards to the user as shown in Figure 4, who can then select an action by simply clicking on the 'Execute' button. The deviation handling action *retry* entails that the same action step has to be retried, e.g., if in Step 2 a deviation occurs and the user selects *retry*, then Step 2 will be executed again. The deviation handling action *give-up* communicates to the robot that assembly process should be stopped and the robot should move to an initialized position. A total of 25 executions of the assembly process by 5 users (roboticists) with 5 executions per user were carried out. Figure 6 shows five different profiles created during online execution. Each profile has three entries, which consists the chosen deviation handling action by that user for Step 2, 3 and 4 respectively. During the online execution of the 5 assembly process per user, different deviations possible were introduced at random such that, overall the user teaches a deviation handling action to the robot for each step of the assembly process.

Learning to handle deviation in this fashion has the following advantages: a) the robot is not required to explore all possible deviation handling actions for each step to find a solution. This reduces the amount of time exponentially as shown in Figure 6 depending on the number of deviation handling actions possible at a given step where a deviation occurred b) the robotic system is capable of handling novel situations (situation not previously encountered) as they occur online during the assembly process with the help of the user interaction. Also, each user has their own profile which will enable the robot to interact with different users according to their chosen fashion.

In the case of a new user, the robot learns deviation handling actions specific to that user and hence collaborates more closely using the user specific preference.



**Fig. 6.** (a) Shows the time taken for the robot to handle deviations in presence and absence of the user assistance. (b) Deviation Handling Profiles of 5 users, learned over a course of 5 assembly process executions per user (total 25 Executions). The three entries in each profile column describe the deviation handling action chosen by that user for the assembly process steps (as given in Figure 5) 2, 3 and 4 respectively.

## 6 Conclusion and Future Work

In this paper we have extended the IRL to enable the robotic system to handle deviations in an assembly process that occur during real-time execution. The eIRL exploits the presence of the human user in the human robot collaboration scenario by interacting with the user for assistance and feedback. The robotic system proposes a set of possible solutions to the user, given a deviation at a particular step in the assembly process. The proposed set of solutions is derived from the knowledge of the assembly process, already known optimal policy, robot capabilities and the set of objects currently present in the workspace. The user assistance is divided into a direct deviation action policy choice and a reward feedback. This enables the system to keep track of user preferences and interact with them more intuitively. An interesting future work would be to combine the already existing preferences of users (the choice of the existing users) and use it to interact with a new untrained user. Also, the evaluation was carried out with roboticists and we would like to extend this evaluation by including users with varying levels of experience with robots.

## 7 Acknowledgment

This research is funded by the projects KoMoProd (Austrian Ministry for Transport, Innovation and Technology), and CompleteMe (FFG, 849441).

## References

1. Akkaladevi, S. C., Heindl, C. Action recognition for human robot interaction in industrial applications. In CGVIS 2015 (pp. 94-99). IEEE (2015)
2. Akkaladevi, S., et. al., Tracking multiple rigid symmetric and non-symmetric objects in real-time using depth data. In ICRA 2016, pp. 5644-5649, (2016)
3. Argall, B., et.al.: A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, vol. 57, no. 5, pp. 469-483 (2009)
4. Bauer, A., Wollherr, D., Buss, M.: Human-robot collaboration: a survey. In: *Int.Journal of Humanoid Robotics*, vol. 5, pp. 47-66 (2008)
5. Dautenhahn, K.: Methodology and themes of human-robot interaction: a growing research field. In: *Int. Journal of Adv. Robotic Sys.* 4, pp. 103-108 (2007)
6. euRobotics.: *Robitics 2020 Strategic Research Agenda for Robotics in Europe*, (2013)
7. Goodrich, M.A., Schultz, A.C.: Human-robot interaction: a survey. In: *Foundations and trends in human-computer interaction*, vol. 1, no. 3, pp. 203-275 (2007)
8. Griffith, S., et.al.: Policy shaping: Integrating human feedback with reinforcement learning. In: *Adv. in Neural Information Processing Systems*, pp. 2625-2633 (2013)
9. Kartoun, U., et.al.: A human-robot collaborative reinforcement learning algorithm. *Journal of Intelligent and Robotic Systems*, vol. 60, no.2, pp. 217-239 (2010)
10. Knox, W. B., et.al.: Interactively shaping agents via human reinforcement: The TAMER framework. In *Proc. Int. Conf. on Knowledge capture*, pp. 9-16 (2009)
11. Knox, W. B., et.al.: Training a robot via human feedback: A case study. In *Social Robotics*, pp. 460-470 Springer (2013)
12. Lee, S., et.al.: Human Mental Models of Humanoid Robots. In: *IEEE Int. Conf. Robotics and Automation*, pp. 2767-2772 (2005)
13. Peng, B., et.al.: A Need for Speed: Adapting Agent Action Speed to Improve Task Learning from Non-Expert Humans. In: *Autonomous Agents*, pp. 957-965 (2016)
14. Rozo, L., et.al.: Learning collaborative impedance-based robot behaviors. In: *Twenty-Seventh AAAI Conference on Artificial Intelligence*, pp. 1422-1428 (2013)
15. Rozo, L., et.al.: Learning force and position constraints in human-robot cooperative transportation. *Robot and Human Interactive Communication*, pp. 619-624 (2014)
16. Scassellati, B.: Theory of mind for a humanoid robot. In: *Autonomous Robots*, vol.12, no. 1, pp. 13-24 (2002)
17. Shen, Y.: *System für die Mensch-Roboter-Koexistenz in der Fließmontage (Forschungsberichte / IWB 305)*, Munich (2015)
18. Suay, H. B., Chernova, S.: Effect of human guidance and state space size on interactive reinforcement learning. In: *2011 Ro-Man*, pp. 1-6 IEEE (2011)
19. Sutton, R. S., and Barto, A. G.: *Reinforcement learning: An introduction*. Vol. 1, no. 1. Cambridge: MIT press (1998)
20. Tenorth, M., et.al.: KnowRob: A knowledge processing infrastructure for cognition enabled robots. *Int. J. of Robotics Research*, vol. 32, pp. 566-590 (2013)
21. Thomaz, A. L., et.al.: Reinforcement learning with human teachers: Understanding how people want to teach robots. In *ROMAN 2006*, pp. 352-357. IEEE (2006)
22. Universal Robots, UR10 robot - A collaborative industrial robot. [Online]: <http://www.universal-robots.com/products/ur10-robot/>, last accessed: 09.08.2016
23. Watkins, C. J., and Dayan, P. Q-learning. *Machine learning*, vol. 8, no. 3-4, pp.279-292.(1992)