

# Active Reward Learning

Christian Daniel\*, Malte Viering\*, Jan Metz\*, Oliver Kroemer\*, Jan Peters\*<sup>†</sup>

\*Institut für Intelligente Autonome Systeme, Technische Universität Darmstadt, 64289, Germany

<sup>†</sup> Max-Planck-Institut für Intelligente Systeme, Tübingen, 72076, Germany

Email: {daniel, kroemer, peters}@ias.tu-darmstadt.de

**Abstract**—While reward functions are an essential component of many robot learning methods, defining such functions remains a hard problem in many practical applications. For tasks such as grasping, there are no reliable success measures available. Defining reward functions by hand requires extensive task knowledge and often leads to undesired emergent behavior. Instead, we propose to learn the reward function through active learning, querying human expert knowledge for a subset of the agent’s rollouts. We introduce a framework, wherein a traditional learning algorithm interplays with the reward learning component, such that the evolution of the action learner guides the queries of the reward learner. We demonstrate results of our method on a robot grasping task and show that the learned reward function generalizes to a similar task.

## I. INTRODUCTION

An important goal of Reinforcement Learning (RL) is to yield more autonomous robots. However, RL methods require reward functions to guide the agent towards a desired behavior. Such reward functions are usually hand coded and, unfortunately, defining rewards for real robot tasks manually is challenging even for relatively well-understood problems such as grasping. Thus, hand coding reward functions is shifting the problem of requiring an expert to design a hard coded controller to requiring a hard coded reward function.

Despite the variety of grasp stability measures that has been developed [34], it has been shown that the resulting grasps are outperformed by grasps learned from kinesthetic teach-in [3, 23]. This example shows that even experts will often design reward functions that are not effective in practice, or lead to undesired emergent behavior [27].

To avoid specifying reward functions, Inverse RL (IRL) extracts a reward function from demonstrations [30, 31, 36]. For many tasks, such demonstrations can be obtained through methods such as kinesthetic teach-in or tele-operation. Unfortunately some skills, such as dynamic skills, are often hard to demonstrate (for example, teaching a robot to throw a basketball). Moreover, both methods require sufficient proficiency of the demonstrator which may lead to the demonstrator having to truly become an expert in performing the task first.

While it may be difficult to analytically design a reward function or to give demonstrations, it is often easy for an expert to rate an agent’s executions of a task. Thus, a promising alternative is to use the human not as an expert in performing the task, but as an expert in evaluating task executions.

Based on this insight, preference based algorithms allow the expert to rank executions and learn a controller based on these rankings [2, 6]. The generally used approach in ranking is to let

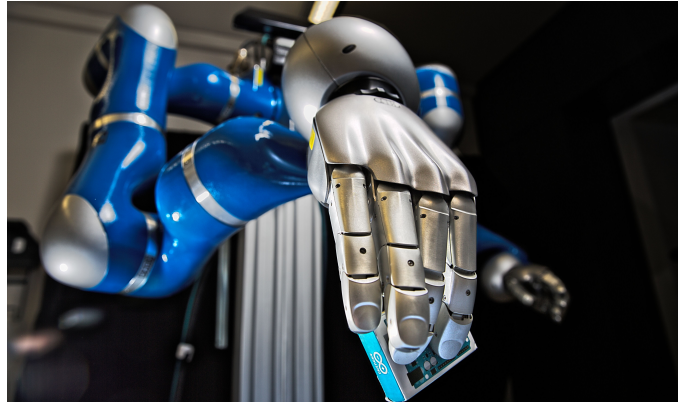


Fig. 1: The Robot-Grasping Task. While grasping is one of the most researched robotic tasks, finding a good reward function still proves difficult.

the expert rank the previously best sample against the current sample, which is an intriguing idea, as humans are better at giving relative judgments than absolute judgments. However, this approach only provides a single bit of information. The technical term describing how much information human subjects can transmit for a given input stimuli is called channel capacity [25]. The channel capacity is a measure for how many input stimuli subjects can distinguish between. In a review of several experiments, Miller [25] concludes that humans have a general channel capacity somewhere between 1.6 and 3.9 bits for unidimensional stimuli. Adding more dimensions to the stimuli further increases this value.

Experiments have also been performed to find out whether humans can transmit more information when labelling stimuli according to prescribed categories or when being allowed to rate on a scale [15]. While there was no significant difference, subjects performed better when rating on a scale.

Based on these insights, we propose an alternative framework in which the human expert can assign numerical values to observed demonstrations. Furthermore, numerical rewards allow to indicate strong preferences over demonstrations. Unfortunately, as already shown by Thomaz and Breazeal [35] and by Cakmak and Thomaz [5], humans have considerable noise in their ratings of actions. To deal with this mismatch between good reward functions for policy search (PS) methods and noisy human rating systems, we propose to learn a probabilistic model of the reward function. In this paper, we take advantage of the Gaussian Process (GP) regression framework as well as the Bayesian Optimization (BO) literature. Having access to BO methods, allows us to efficiently minimize

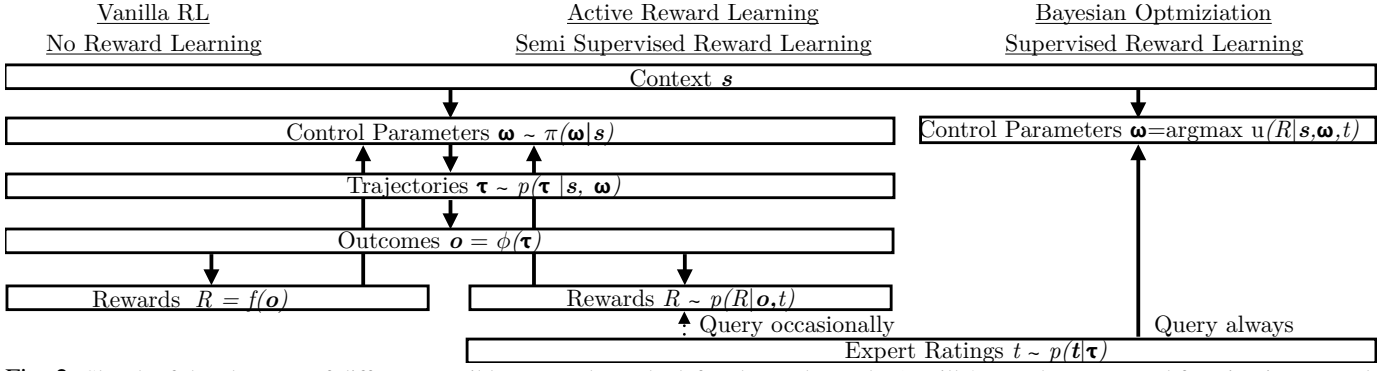


Fig. 2: Sketch of the elements of different possible approaches. The left column shows the ‘vanilla’ RL where a reward function is assumed. The middle column shows the proposed active reward learning approach, which shares the policy learning component with vanilla RL but models a probabilistic reward model that gets updated by asking for expert ratings sometimes. The right column shows the BO approach, where actions are chosen that maximize the utility function (for example one of the acquisition functions presented in II-A1 and requires an expert rating for each of the chosen actions).

the number of expert interactions, which is essential when designing methods for ‘autonomous’ agents. We evaluate the proposed method on a series of simulated and one real robot task and we also show that a reward function which is learned for one task (grasping a box) can be used to learn a new task (grasping a pestle).

## II. METHOD

While there are many approaches to RL, Policy Search (PS) is especially well suited for real robot tasks. It does not rely on exhaustive sampling of the state-action space and, as a result, many recent advances in robot learning rely on policy search [22, 20, 26]. We consider the contextual episodic PS case with continuous contexts  $s \in \mathcal{S}$  and continuous control parameters  $\omega \in \Omega$ . We use the term ‘context’ instead of the more general term ‘state’ to indicate the initial state of the robot and the environment. The distribution over contexts is denoted by  $\mu^\pi(\cdot)$ . Whereas the traditional RL notation uses actions  $\mathbf{a}$ , we instead write control parameters  $\omega$  to clarify that we directly optimize for the parameters of a controller which can, for example, be a Movement Primitive as discussed in Section III-B. Executing a rollout with control parameters  $\omega$  results in a trajectory  $\tau \sim p(\tau|s, \omega)$ , where the trajectory encodes both, the robot’s state transitions as well as relevant environment state transitions. The agent starts with an initial control policy  $\pi_0(\omega|s)$  in iteration 0 and performs a predetermined number of rollouts. After each iteration, the agent updates its policy to maximize the expected reward

$$\mathbb{E}_{s, \omega, \tau}[R(\tau)] = \iiint R(\tau) p(s, \omega, \tau) d\tau ds d\omega, \quad (1)$$

with  $p(s, \omega, \tau) = p(\tau|s, \omega)\pi(\omega|s)\mu^\pi(s)$ . Parameter-based PS has been shown to be well suited for complex robot tasks [9], as it abstracts complexity while retaining sufficient flexibility in combination with suitable movement primitive representations.

In the original RL setting a reward function is assumed to be known that evaluates the agents behavior and assigns rewards to the agent’s actions. While not always explicitly stated, the reward function for real robot tasks often depends not on the whole trajectory  $\tau$  of the robot, but on features  $\phi(\tau)$  of the trajectory. We refer to the result of the feature extraction from

trajectories as the *outcomes*  $\mathbf{o} = \phi(\tau)$ . We assume that the reward depends only on the features of the trajectories. Such outcomes could, for example, describe the minimal distance to goal positions or the accumulative energy consumption.

As the problem of specifying good features  $\phi(\tau)$  is beyond the scope of this paper, we assume the features to be known, as is usually assumed when designing reward functions.

As outcomes are of much lower dimensionality than trajectories, we can use outcomes to efficiently model the reward function  $\hat{R}(\mathbf{o})$  from a training set  $\mathcal{D} = \{\mathbf{o}_{1:n}, \mathbf{R}_{1:n}\}$  using regression. Using the feature representation, the term  $R(\tau)$  in Eq. 1 is replaced by  $R(\mathbf{o} = \phi(\tau))$  to become

$$\mathbb{E}_{s, \omega, \tau}[R(\mathbf{o})] = \iiint R(\mathbf{o} = \phi(\tau)) p(s, \omega, \tau) d\tau ds d\omega. \quad (2)$$

Obviously, it is not sufficient to build the reward function model on samples from the initial policy  $\pi_0(\omega|s)$ , as the agent is likely to exhibit poor performance in the early stages of learning and the reward learner would not observe good samples. Therefore, we need to couple the process of learning a good policy  $\pi(\omega|s)$  and a good reward function  $\hat{R}(\mathbf{o})$ , such that they are developed interdependently. However, in such a coupled active learning process, we need to know which samples to query from our expert, as we want to minimize expert interaction.

Modelling a probabilistic reward model  $p(R|\mathbf{o}, \mathcal{D})$ , instead of a deterministic reward function  $R(\mathbf{o})$ , allows us to leverage the information about the certainty of our estimate to control the amount of human interaction required, as we only need to interact with the human expert if our model of the reward is not certain enough. We describe the details of this process in Section II-A.

When using a probabilistic model of the reward, we have to replace the reward term in Eq. 1 with

$$R(\mathbf{o}) = \mathbb{E}_{p(R|\mathbf{o})}[R] = \int p(R|\mathbf{o}) R dR,$$

as most PS methods work on the expectation of the reward. Using the probabilistic model of the reward function the PS method continuously learns how to achieve high reward outcomes, while the reward model learner adapts the accuracy of its model.

<b>Input:</b> Information loss tolerance $\epsilon$ , improvement threshold $\lambda$ , acquisition function $u$
<b>Initialize</b> $\pi$ using a single Gaussian with random mean. GP with zero mean prior.
<b>while</b> not converged
<b>Set sample policy:</b> $q(\omega \mathbf{s}) = \pi_{\text{old}}(\omega \mathbf{s})$
<b>Sample:</b> collect samples from the sample policy $\{\mathbf{s}_i \sim p(\mathbf{s}), \omega_i \sim q(\omega \mathbf{s}_i), \mathbf{o}_i\}, i \in \{1, \dots, N\}$
<b>Define Bellman Error Function</b> $\delta(\mathbf{s}, \omega) = R(\mathbf{o}) - V(\mathbf{s})$
<b>Minimize the dual function</b> $[\alpha^*, \eta^*] = \arg \min_{[\alpha, \eta]} g(\alpha, \eta)$
<b>Determine base line</b> $V(\mathbf{s}) = \alpha^{T*} \phi(\mathbf{s})$
<b>Policy update:</b> Calculate weighting $p(\mathbf{s}_i, \omega_i) \propto q(\mathbf{s}_i, \omega_i) \exp\left(\frac{1}{\eta^*} \delta^*(\mathbf{s}_i, \omega_i)\right)$ Estimate distribution $\pi(\omega \mathbf{s})$ by weighted maximum likelihood estimates
<b>Reward model update:</b> FindNominee = true <b>while</b> FindNominee
<b>Nominate outcome:</b> $\mathbf{o}^+ = \arg \max u(\mathbf{o})$ <b>if</b> $(\mathbf{o}^+ \notin \mathcal{D}) \wedge (\sigma(\mathbf{o}^+)/\beta > \lambda)$ Demonstrate Corresponding Trajectory $\tau^+$ Query Expert Reward $R^+$ $\mathcal{D} = \mathcal{D} \cup \{\mathbf{o}^+, R^+\}$ <b>else</b> FindNominee = false
Update reward model $p(R \mathbf{o}, \mathcal{D})$ Optimize GP-hyper parameters $\theta$
<b>Output:</b> Policy $\pi(\omega \mathbf{s})$ , reward model $p(R \mathbf{o}, \mathcal{D})$

TABLE I: We show the algorithmic form of active reward learning with REPS. We specify the information loss tolerance  $\epsilon$  as well as an initial sampling policy and an improvement threshold  $\lambda$ . In each iteration, the algorithm first samples from the sampling policy, minimizes the dual function to find values for  $\alpha^{T*}$  and  $\eta^*$  and then computes the next policy. After each policy search iteration, the reward function learner chooses whether to demonstrate samples to the expert according to the acquisition function. The parameters  $\alpha$  and  $\eta$  are parameters of the dual function problem of REPS and can be optimized through standard optimization algorithms [9].

### A. Active Reward Learning

Our goal is to find a model  $p(R|\mathbf{o}, \mathcal{D})$  that predicts the reward given an observed outcome and training data  $\mathcal{D}$ , which is obtained from an expert. When modelling the reward, we have to take into account that the expert can only give noisy samples of his implicit reward function and we also have to model this observation noise. Thus, we need to solve the regression problem

$$R(\mathbf{o}) = f(\mathbf{o}) + \eta, \quad \eta \sim \mathcal{N}(0, \beta),$$

where we assume zero mean Gaussian noise. Such a regression problem can, for example, be solved with Gaussian Process (GP) regression

$$f(\mathbf{o}) \sim \mathcal{GP}(m(\mathbf{o}), k(\mathbf{o}, \mathbf{o}')),$$

where  $m(\mathbf{o})$  is the mean function and  $k(\mathbf{o}, \mathbf{o}')$  is the covariance function of the GP. For the remainder of this paper we

use the standard squared exponential covariance function

$$k(\mathbf{o}, \mathbf{o}') = \theta_0^2 \exp\left(-\frac{\|\mathbf{o} - \mathbf{o}'\|^2}{2\theta_1^2}\right).$$

The set of hyper parameters  $\theta = \{\theta_0, \theta_1, \beta\}$  is found through optimization [29]. Given a training set  $\mathcal{D} = \{\mathbf{o}_{1:n}, \mathbf{R}_{1:n}\}$ , we can write down the covariance matrix between previously observed rewards and outcomes

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{o}_1, \mathbf{o}_1) & \dots & k(\mathbf{o}_1, \mathbf{o}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{o}_n, \mathbf{o}_1) & \dots & k(\mathbf{o}_n, \mathbf{o}_n) \end{bmatrix} + \beta \mathbf{I}.$$

Assuming a zero mean prior, the joint Gaussian probability of the training samples in  $\mathcal{D}$  and the reward prediction  $R^+$  of a new unrated observation is given by

$$\begin{bmatrix} \mathbf{R}_{1:n} \\ R^+ \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \hat{\mathbf{k}} \\ \bar{\mathbf{k}} & k(\mathbf{o}^+, \mathbf{o}^+) \end{bmatrix}\right),$$

with  $\hat{\mathbf{k}} = [k(\mathbf{o}_1, \mathbf{o}^+) \dots k(\mathbf{o}_n, \mathbf{o}^+)]^T$  and  $\bar{\mathbf{k}} = [k(\mathbf{o}^+, \mathbf{o}_1) \dots k(\mathbf{o}^+, \mathbf{o}_n)]$ . The predictive posterior reward  $p(R^+|\mathbf{o}, \mathcal{D})$  of a new outcome  $\mathbf{o}^+$  is then given by a Gaussian

$$p(R^+|\mathbf{o}, \mathcal{D}) \sim \mathcal{N}(\mu(\mathbf{o}^+), \sigma^2(\mathbf{o}^+)),$$

with mean and variance

$$\begin{aligned} \mu(\mathbf{o}^+) &= \mathbf{k}^T \mathbf{K}^{-1} \mathbf{R}_{1:n}, \\ \sigma^2(\mathbf{o}^+) &= k(\mathbf{o}^+, \mathbf{o}^+) - \mathbf{k}^T \mathbf{K}^{-1} \mathbf{k}, \end{aligned}$$

by conditioning the GP on the observed outcome  $\mathbf{o}^+$ . Using the GP, we can represent both our expected reward  $\mu(\mathbf{o}^+)$ , which is provided to the policy learner, and the variance of the reward  $\sigma^2(\mathbf{o}^+)$  which is essential to the active learning component. The reward variance  $\sigma^2(\mathbf{o}^+)$  depends on the distance of the outcome  $\mathbf{o}^+$  to all outcomes in the training set  $\mathcal{D}$  and the observation noise  $\beta$ , which is a hyper parameter that we optimize for. Using the predictive variance, we can employ one of many readily available BO methods to find the maximum of the reward function.

1) *Optimizing the Reward Model:* The goal of BO is to optimize a function under uncertainty. Acquisition functions (AFs), are utility functions which maximize their function value at locations of the input space which are likely to maximize the original problem. AFs usually encode an exploration-exploitation trade-off, such that they do not only query samples in known high value regions but also in regions that have not been sufficiently explored before. While using GPs to model the reward function allows us access to the BO toolbox in general, we deviate from the standard BO approach in two points. First, as our GP models a relationship of outcomes to rewards instead of context-actions to rewards, we cannot sample arbitrary outcomes  $\hat{\mathbf{o}}$  to improve our estimate. To do so, we would require access to  $p(\tau|\mathbf{s}, \omega)$  such that we can request the agent to perform actions that result in trajectories  $\hat{\tau}$  which yield the outcome  $\hat{\mathbf{o}} = \phi(\hat{\tau})$ . Second, we would need to guarantee that the outcomes requested by the AF to improve the reward model are physically possible. However,

this transition model is unknown and, thus, we have to revert to the previously observed outcomes that have been generated during the agent’s learning process so far.

Furthermore, we need to balance the improvement of our current estimate of the reward function and the number of queries that we request from the expert, i.e., we want to find a trade-off between finding a good reward function and learning the task with minimal human input.

2) *Sample Efficiency*: In the traditional BO framework, the goal is to find a global maximum. Sampling of the function can be stopped when the improvement of the function is marginal, for example when the predicted variance around the optimum is very low. However, in the proposed scenario where a policy  $\pi(\omega|s)$  and a reward model  $p(R|o)$  are learned simultaneously and obtaining training samples of the reward model is very expensive, the problem of deciding *when* to improve the reward model becomes crucial.

The reward model relies on the policy  $\pi(\omega|s)$  to provide outcomes in interesting, i.e., high reward regions, and the policy relies on the reward model  $p(R|o)$  to guide its exploration towards such regions of interest. Thus, the reward model needs sufficient training data to approximately predict the reward function in early stages and a higher density of training points once the agents policy starts to converge to a solution. At the same time, we want to minimize the number of queries to the expert over the learning process.

In order to balance this trade-off, we propose an acquisition algorithm which, according to a selected acquisition function  $u(o)$ , first samples the best observed sample outcome from the history of the agent’s outcomes  $o_{n+1} = \arg \max_o u(o|\mathcal{D})$ . In this sample based search, we also include all outcomes that have been demonstrated to the expert. If the acquisition function is maximized by an already demonstrated outcome, we stop and do not query any samples in this iteration. Maximizing the AF by an already observed outcome is unique to the sample based case, as in the continuous case, a point around an observed outcome would usually have more variance and, thus, maximize the AF. With our approach we achieve a sparse sampling behavior that requires less expert interactions. If, however, the outcome that maximizes  $u$  has not yet been rated by the expert, we need to decide whether querying the outcome is beneficial. For example, we may have already converged to a good estimate of the reward function, and new outcomes improve on the mean reward solely due to the observation noise. In this case we do not want to query the expert. Thus, we decide whether to query an outcome by thresholding the ratio of predictive variance and estimated observation noise  $\sigma(o)/\sqrt{\beta} > \lambda$ , where  $\lambda$  is a tuning parameter which allows us to trade off the accuracy of the final solution with the query frequency by adapting the available AFs to explicitly take the estimated observation noise into account. While this adaptation of the AFs introduces a new parameter to be tuned, our experiments show that the use of this technique results in less human interactions while maintaining high performance and tuning of the parameter is straightforward.

If we decide to query the sample, we update the GP and

search for the new maximum of the acquisition function, otherwise we stop and do not update the GP further in this iteration.

The general information flow of our method is as follows. We start with an uninformed, i.e., zero mean GP for  $p(R|o)$  and we initialize the PS method with an initial policy  $\pi_0(\omega|s)$ . The PS learner then starts performing one iteration of rollouts, which we also call episodes. After each iteration of rollouts, rewards for the outcomes of the resulting trajectories  $o = \phi(\tau)$  are requested from the reward learner  $R \sim p(R|o)$ . The reward learner then decides whether to ask for expert ratings for any of the outcomes to update its model. In that case, the agent repeats the corresponding episode to present the outcome to the expert. Finally, the reward learner returns the mean estimate of the reward to the PS learner, which uses the rewards to update its policy and start the next iteration.

### III. BACKGROUND

In this section we provide compact background information on components of the proposed algorithms that are necessary to give a complete picture of the proposed approach.

#### A. Relative Entropy Policy Search

We pair our reward learning algorithm with the recently proposed relative entropy policy search (REPS) [28]. REPS is a natural choice for a PS method to be combined with an active learning component as it has been shown to work well on real robot problems [9] and is designed to ‘stay close to the data’. Thus, previous expert queries will remain informative. A distinctive feature of Relative Entropy Policy Search (REPS), is that its successive policies vary smoothly and do not jump in the parameter space or the context space. This behavior is a beneficial characteristic to increase compliance with our proposed active learning approach, as we are only able to predict correct rewards within observed regions of the parameter space. To constrain the change in the policy, REPS limits the Kullback-Leibler divergence between a sample distribution  $q(s, \omega)$  and the next distribution  $\pi(\omega|s)\mu^\pi(s)$

$$\epsilon \geq \sum_{s, \omega} \mu^\pi(s) \pi(\omega|s) \log \frac{\mu^\pi(s) \pi(\omega|s)}{q(s, \omega)}. \quad (3)$$

For the complete optimization problem and its solution we refer to the original work [28].

#### B. Dynamic Movement Primitives

A popular use case for PS methods is to find good parameters of trajectory generators such as Dynamic Movement Primitives (DMPs) [17]. The resulting desired trajectories can then be tracked by a linear feedback controller. DMPs model trajectories using an exponentially decreasing phase function and a nonlinear forcing function. The forcing function excites a spring damper system that depends on the phase and is guaranteed to reach a desired goal position, which is one of the parameters of a DMP. The forcing function is modelled through a set of weighted basis functions  $\omega\Psi$ . Using the weights  $\omega$  of the basis functions  $\Psi$  as parameters, we can

learn parametrized joint trajectories, and an increasing number of basis functions results in an increased flexibility of the trajectory. We use DMPs for our simulated robot experiments.

### C. Bayesian Optimization for RL

An alternative approach to learning control parameters  $\omega$  using PS methods is to model  $p(R|s, \omega)$  directly and to use Bayesian Optimization (BO) instead of the PS learner. However, the approach proposed in this paper introduces a layer of abstraction which allows us to learn a mapping of only a low-dimensional input space to the reward, as we only have to map from outcomes to reward as opposed to map from state-action to reward. As BO methods are global methods, they are more susceptible to the curse of dimensionality than PS methods which are local methods. As a result, the mapping  $p(R|s, \omega)$  which the standard BO solution uses is considerably more difficult to learn than the mapping of the modular approach that we are proposing.

The BO approach would also require expert ratings for every sample, while the proposed approach requires only occasional human feedback.

### D. Acquisition Functions

In the following we present four Acquisition Function (AF) schemes taken from Hoffman et al. [16] that we used to optimize the model of the reward function.

a) *Probability of Improvement*: The Probability of Improvement (PI) [16] in its original formulation greedily searches for the optimal value of the input parameter that maximizes the function. An adapted version of PI balances the greedy optimization with an exploration-exploitation trade-off parameter  $\xi$ . The adapted version is given by

$$\text{PI}(\mathbf{o}) = \Phi\left(\frac{\mu(\mathbf{o}) - f(\mathbf{o}^*) - \xi}{\sigma(\mathbf{o})}\right),$$

where  $\mathbf{o}^*$  is the best sample in the training set  $\mathcal{D}$  and  $\Phi(\cdot)$  is the normal cumulative density function. The exploration-exploitation trade-off parameter  $\xi$  has to be chosen manually.

b) *Expected Improvement*: Instead of finding a point that maximizes the PI, the Expected Improvement (EI) [16], tries to find a point that maximizes the magnitude of improvement. Thus, it does not only try to improve local maxima but also considers maxima in different regions and is less greedy in the search of an optimal reward  $R$ .

$$\text{EI}(\mathbf{o}) = (\mu(\mathbf{o}) - f(\mathbf{o}^*) - \xi) \Phi(M) + \sigma(\mathbf{o})\rho(M),$$

if  $\sigma(\mathbf{o}) > 0$  and zero otherwise, where  $\rho(\cdot)$  is the normal probability density function.  $M$  is given by

$$M = \frac{\mu(\mathbf{o}) - f(\mathbf{o}^*) - \xi}{\sigma(\mathbf{o})}.$$

The EI acquisition function shares the tuning factor  $\xi$  for the exploitation-exploration trade-off with the PI, where a suggested value is  $\xi = 0.01$  [16].

c) *Upper Confidence Bound*: The Upper Confidence Bound directly uses the mean and standard deviation of the reward function at the sample location to define the acquisition function. An adapted version of the UCB function [33] is given by

$$\text{GP-UCB}(\mathbf{o}) = \mu(\mathbf{o}) + \sqrt{v\gamma_n\sigma(\mathbf{o})},$$

where recommended values  $v = 1$  and  $\gamma_n = 2\log(n^{d/2+2}\pi^2/3\delta)$  with  $d = \dim(\mathbf{o})$  and  $\delta \in (0, 1)$  are given by Srinivas et al. [33].

d) *GP Hedge*: As each of the above acquisition functions lead to a characteristic and distinct sampling behavior, it is often not clear which acquisition function should be used. Portfolio methods, such as the GP-Hedge [16], evaluate several acquisition functions before deciding for a sample location. Given a portfolio with  $J$  different acquisition functions, the probability of selecting acquisition function  $j$  for the sample  $n + 1$  is given by the softmax

$$p(j) = \frac{\exp(\eta g_n^j)}{\sum_{i=1}^J \exp(\eta g_n^i)},$$

where  $\eta > 0$  is the temperature of the soft-max distribution. The gains vector  $\mathbf{g}$  is initialized to zero before taking the first sample  $g_0^{1:J} = 0$ , and is then updated with the cumulative reward gained by the selected acquisition function, i.e.,  $g_{n+1}^j = g_n^j + \mu(\mathbf{o}_{n+1}^j)$ , where  $\mathbf{o}_{n+1}^j$  is the sample point nominated by acquisition function  $j$ . For all other gains the value does not change, i.e.,  $g_{n+1}^{i \neq j} = g_n^i$ .

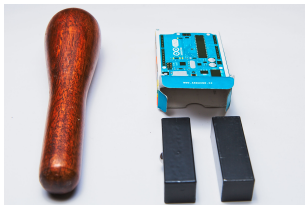
## IV. EVALUATIONS

In this section we show evaluations of the proposed active reward learning approach. For all simulation experiments, unless stated otherwise, we tested each setting ten times.

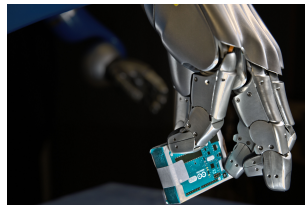
### A. Five Link Reaching Task

A simulated planar robot consisting of five links connected by rotary joints was controlled in joint space using Dynamic Movement Primitives (DMP) [17], as described in Section III-B. If not stated otherwise, we used 20 basis functions per joint, resulting in a total of 100 parameters that had to be learned. We evaluated our approach on a reaching task, where an analytical reward function was readily available. The hand coded reward function was given by  $R(\mathbf{p}_r) = 1000 - 100\|\mathbf{p}_r - \mathbf{p}_g\|$ , where  $\mathbf{p}_r$  was the position of the robot's end effector and  $\mathbf{p}_g$  was the desired target position. However, we increased the difficulty of the task by not supplying the outcome features in task space, but rather in joint space, i.e., the GP had to model the forward kinematics to predict the reward, making the problem both non-convex and high-dimensional (five dimensional mapping).

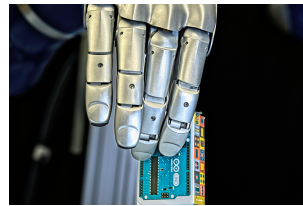
To allow extensive and consistent evaluation of all parameters of the presented approach, we programmed a noisy expert, which returned the reward with additive white noise (standard deviation was 20). In Fig. 7 we present results that compare the coded noisy expert approach to actually querying a human expert and show that the behavior is comparable. The human



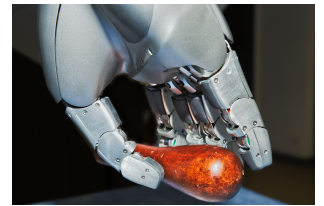
(a) Pestle and the paper box that is filled with metal bars.



(b) Failed grasp, not robust against perturbations.



(c) Mediocre grasp, stable but incorrect orientation.



(d) Good grasp, stable with intended orientation.

Fig. 3: Examples of different grasps and their categorization. Grasps count as failed if the object is either not picked up at all or if small perturbations would make the object drop. Grasps that are stable but do not keep the original orientation of the object count as OK but not successful grasps. Grasps that are both stable and keep the original orientation count as successful grasps.

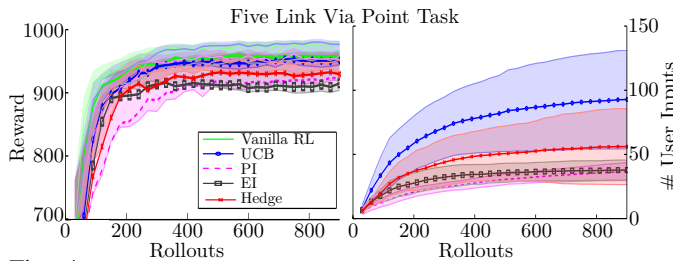


Fig. 4: We evaluated our approach on a programmed, but noisy expert to emulate human expert input. Vanilla RL REPS queries the programmed expert for every sample, while our approach builds a model of the reward function and only queries the expert occasionally.

expert could give rewards on a vertical bar through a graphical interface.

1) *Evaluation of Acquisition Functions:* The evaluation of the available AFs given in Fig. 4 shows that even though there was only limited difference in the asymptotic performance, there was considerable difference in the sample efficiency. Especially the GP-UCB AF asks for many user queries. While PI had the lowest asymptotic performance, as it is the most greedy of the presented AFs, it also required the lowest number of user queries. This behavior makes it an interesting candidate when trying to minimize human interactions.

2) *Evaluation of Uncertainty Threshold:* To optimally trade off the number of queries and the agents performance we need to set the uncertainty threshold trade-off parameter  $\lambda < \sigma(o)/\sqrt{\beta}$ . This parameter expresses how certain we require our algorithm to be that a proposed query is not explained by the estimated observation noise. If we choose  $\lambda$  to be greater than 1, we require our estimate to improve on the observation noise. Fig. 5 show the effects of adjusting  $\lambda$ . The performance remained stable up to  $\lambda = 1.3$  and started degrading with  $\lambda = 2$ . Changing the order of magnitude of  $\lambda$  resulted in a failure to learn the task. While the effects of setting  $\lambda = 1.5$  on the performance were moderate, the number of queries were reduced by about 50% when compared to  $\lambda = 1.3$ .

3) *Evaluation of Sparse Sampling:* In order to minimize human interaction, we stop improving the GP in each iteration if the outcome that maximizes the AF has already been queried before, instead of selecting the second best outcome according to the AF. The results in Fig. 6 show, that sparse sampling leads to equally good asymptotic performance but requires considerably less expert interactions.

4) *Evaluation Direct Learning:* The premise of this paper is that while BO can be used to efficiently learn the mapping

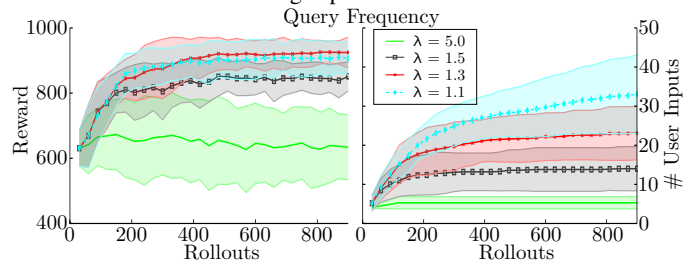


Fig. 5: We evaluated the impact of the threshold factor  $\lambda$  that influences when and how often we sample. Lower values of  $\lambda$  led to better converged performance but required more user interaction.

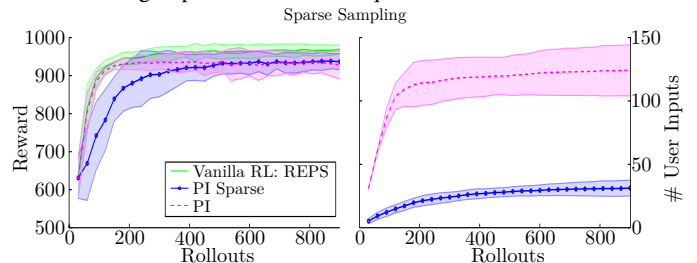


Fig. 6: Results of comparing our acquisition algorithm to the standard acquisition algorithm. Our algorithm collects sparse user queries and does not ask for any user queries after a PS iteration if the outcome that maximizes the AF has already been queried before.

of outcome to reward, it would be too sample intensive to directly learn the mapping from parameter space to outcome. We validated this premise by comparing our joint approach to directly learning the reward from the control parameters (i.e., from the basis function weights  $\omega$ ). The results in Fig. 8 show that BO did not converge to a good solution within 50 expert queries. Both methods were using PI.

5) *Comparison to Inverse Reinforcement Learning:* We compared our algorithm to the Maximum Entropy IRL approach [36] on a via point task in two different scenarios. Trajectories generated by DMPs had to pass one or two via points, respectively (20 dimensional action space for each).

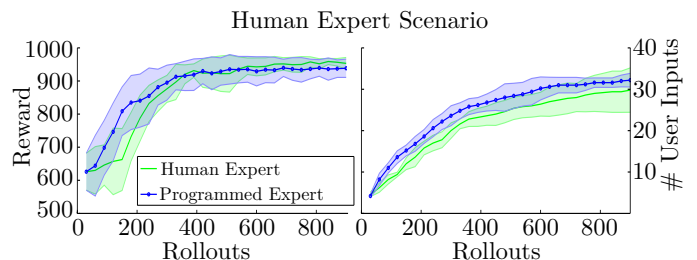


Fig. 7: We validated our approach of using a noisy programmed expert as substitute for a human expert on the simulated tasks. The results show that both experts yielded similar learning behavior.



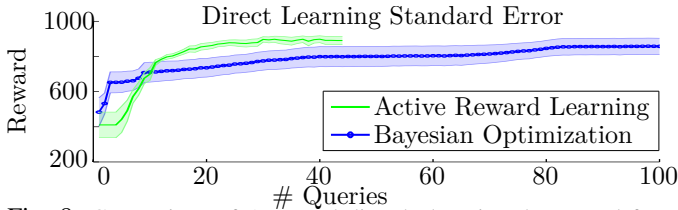


Fig. 8: Comparison of ARL and directly learning the reward from the control parameter  $\omega$  using BO. In this figure we plot the standard error as opposed to the standard deviation. ARL performs significantly better after 50 queries ( $p < 0.05$ ).

		#UI = 5	#UI = 10	#UI = 20
ARL	1VP	<b>995.2 ± 2.477</b>	<b>999.2 ± 0.25</b>	<b>999.2 ± 0.25</b>
IRL	1VP	983 ± 3.306	977.3 ± 6.334	984 ± 3.036
ARL	2VPs	970.7 ± 13.1	<b>998 ± 1.13</b>	<b>999.1 ± 0.172</b>
IRL	2VPs	<b>994.1 ± 1.874</b>	996 ± 0.4736	996 ± 0.555

TABLE II: Comparison of IRL and ARL. We compare the methods on two via point (VP) tasks with one and two VPs. The columns show the achieved performance (mean and standard deviation) after five, ten or 20 user interactions (UIs).

For this comparison only, we relaxed our assumption that no demonstrations are available and provided the IRL approach with (imperfect) demonstrations that passed close to the via point (at most 0.1 units away). In Table II, we show the mean and standard deviation of the best policy reached for either five, ten or 20 user interactions (UI). In our approach, UIs are ratings while in IRL UIs are demonstrations. The results show that our approach yields competitive results while not requiring access to demonstrations.

6) *Alternative Policy Learners:* While we used REPS for most of our experiments, the proposed framework is not limited to a specific RL method. To investigate the compatibility with other methods, we compared our framework using REPS with our framework in combination with a Bayesian Policy Gradient (BPG) method [13] on a via point task with one via point (20 dimensional action space). BPG models the gradient of the policy using a GP and uses the natural gradient in the policy update. The results in Fig. 9 show that both methods were able to learn the task in combination with our approach.

### B. Robot Grasping

We used the results from Section IV-A to set the parameters for the real robot experiment (we use PI and set  $\lambda = 3$ ). We learned the policy  $\pi(\omega|s)$  with 15 samples per iteration for a total of 10 iterations and we repeated the experiment three times. The control parameters of the policy were the 15 joints of the five finger DLR hand, which is mounted to a 7 DOFs

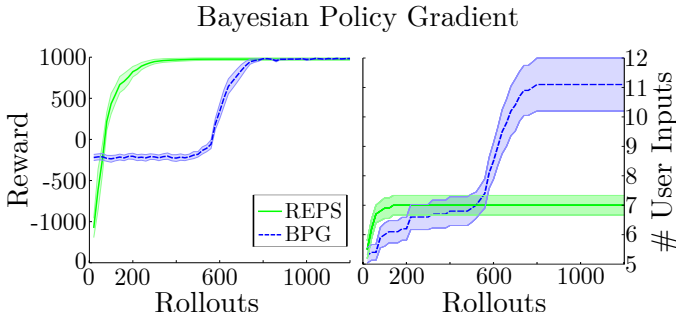


Fig. 9: Performance of REPS and BPG on a via point task.

KUKA lightweight arm as shown in Fig. 1. We considered the task of blind grasping as described in [8], where no object information was available and we did not have visual feedback, i.e., we did not have information about the contact points. Instead, we calculated the forces in the finger tips through the joint torques and the hand kinematics. We used the finger tip force magnitudes as outcome features which were used by the GP to model the reward function.

Evaluating the real robot experiments presented the problem of a success metric. As we did not have a ‘correct’ reward function that we could evaluate the learned reward function against, we resorted to introducing three label categories which were only used for the evaluation after finishing the experiments. The scheme presented in Fig. 3 labels grasps as failures with a reward of  $-1$ , if the object was not lifted at all or slipped when slightly perturbed. Grasps that were stable but did not keep the intended orientation were given a reward of 0. Finally, grasps that lifted the object and kept the orientation of the object were assigned a reward of 1. These labels and reward values were not used during the learning of either the policy or the reward model but only used to present results. During the learning phase, the human expert assigned grasp ratings in the range of  $\pm 1000$ .

1) *Learn to Grasp Unknown Object:* The object to be grasped was a cardboard box filled with metal weights such that the robot cannot grasp the object with very unstable grasps or by deforming the paper box. The box, shown in Fig. 3a, was of size 7.5cm x 5.5cm x 2cm and filled with two metal bars with a combined weight of 350g. The results of three trials presented in Fig. 10 show that the robot learned to perform a successful grasp of the object in all three trials, while only requesting six queries in the first and twelve queries in the second trial. In the last trial, the robot’s performance first increased quickly but dropped after 80 episodes (or rollouts), coinciding with a sudden increase of user queries, such that the final number of queries in the last trials was 27. The reason for this unusual behavior was a malfunction of the distal joint of the thumb, which rendered the grasping scheme dysfunctional. As the learner could not reproduce outcomes that led to good rewards, it resorted to finding a different grasp strategy. At the same time, since the GP was presented with new outcome samples in previously unobserved regions of the outcome space, it requested new user queries to model the reward function in the new region of interest.

We compared our learned reward function to a (naive) hand coded reward function based on the same features. The hand coded reward function aimed to reach a total force magnitude over all fingers. Using the programmed reward, the robot was able to reliably pick up the object after the first two trials and in ten out of 15 grasps at the end of the third trial (We expect the robot would have also learned to pick it up reliably in the last trial with more iterations). However, the robot did not pick up the object in a way that kept the original orientation of the object. Encoding such behavior through only force features by hand is challenging. The performance curve of the hand coded reward function shows a slight dip, which is possible as we

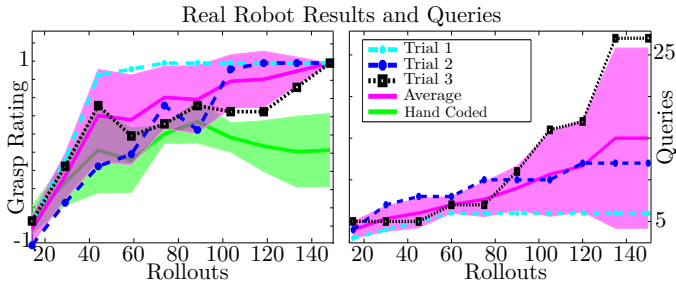


Fig. 10: The real robot successfully learned good grasps in all of the three trials. The grasp rating scheme is described in Fig. 3. In the last trial, one joint failed, represented by a spike in expert queries. The robot successfully adapted the reward function to recover from the hardware failure. We also show the average performance of a naive hand coded reward function.

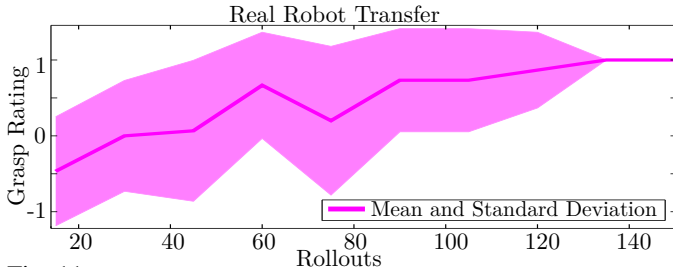


Fig. 11: Performance of one trial on the real robot system. The robot uses the reward function learned in trial one of the original task to learn to grasp a new, unknown object (a wooden pestle as shown in Fig. 3a). For this trial, we did not allow any expert queries. The robot successfully learned to grasp the object.

are not plotting the internal reward but the reward assigned according to the grading scheme introduced in Fig. 3.

2) *Transfer Learned Reward Function to New Object:* As we based our reward model only on the finger tip forces, it is modular and can also be used for (reasonably) different objects. To test these generalization capabilities, we started a new trial with a different object (a pestle as shown in Fig. 3a) and initialized the GP with the training data from the first trial of the previous task. For this experiment we only used ratings from the previous trial, and did not ask for additional expert ratings. The pestle that we used for this task was similar in dimensions but different in shape, such that the agent had to learn different joint configurations to achieve similar finger forces that optimized the reward function. The dimensions of the pestle were 18cm length, 1.5cm radius on the thin end and 2.25cm radius on the thick end. The results in Fig. 11 show that the agent was able to learn to reliably perform a robust grasp on the pestle.

## V. RELATED WORK

The term reward shaping has been used to describe efforts to adapt rewards functions such that the resulting policy remains invariant but learning speed is increased [27, 10, 4, 32] or that learning on a new task is accelerated [21]. Aside from reward shaping, Dorigo and Colombetti [11] have used the term robot shaping to describe efforts of teaching different tasks to a robotic agent. Derived from this terminology the TAMER framework [19] uses the term interactive shaping. In TAMER the agent receives reinforcements from a human, guiding the learning process of the agent, but there is no active

component. The Advise framework [14] also uses the term shaping to describe the process of modelling an oracle and are actively requesting labels for an agent’s action (good or bad). This approach is tailored for discrete settings and the feedback frequency is pre-defined.

In preference learning algorithms, the expert is usually requested to rank the current best against a new execution and the reward function is inferred from these rankings [2, 6]. Akrouf et al. [1] can also deal with noisy rankings. Chu and Ghahramani [7] introduced the use of GPs to preference rankings. In preference based approaches, the expert is limited to transmit one bit of information and cannot express strong preferences, as, supposedly, human experts cannot give numerical rewards that are sufficiently accurate. Our contribution is to show how a rating based approach with explicit noise model can be used in real robot continuous state-action space problems taking advantage of the stronger guidance through strong preferences (large differences in assigned rewards).

The problem of expert noise has also been addressed in Bayesian RL and IRL. In Bayesian RL, Engel et al. [12] has proposed to model the value function through a GP and Ghavamzadeh and Engel [13] has proposed to model the policy gradient through a GP. Especially the policy gradient method is also well suited to work in combination with our approached framework. If demonstrations are available, IRL is a viable alternative. Ziebart et al. [36] have relaxed the assumptions on the optimality of demonstrations such that a reward function can be extracted from noisy demonstrations.

In a combination of IRL and preference learning, Jain et al. [18] have proposed the iterative improvement of trajectories. In their approach, the expert can choose to rank trajectories or to demonstrate a new trajectory that does not have to be optimal but only to improve on the current trajectory. This approach cannot directly be used on learning tasks such as grasping, as a forward model is required. Alternatively, Lopes et al. [24] propose a framework where IRL is combined with active learning such that the agent can decide when and where to ask for demonstrations.

## VI. CONCLUSION & FUTURE WORK

We presented a general framework for actively learning the reward function from human experts while learning the agent’s policy with any PS or policy gradient method. Our experiments showed that the learned reward function outperforms naive hand coded reward functions, generalizes to similar tasks and that the approach is robust to sudden changes in the environment, for example when a mechanical failure occurs. In future work we plan to extend on the ground-laying work of this paper and investigate possible synergies between specific PS methods and the reward learning framework to improve on the performance of the generic AFs.

## ACKNOWLEDGMENTS

The authors want to thank for the support of the European Union projects # FP7-ICT-270327 (Complacs) and # FP7-ICT-2013-10 (3rd Hand).



## REFERENCES

- [1] Riad Akrou, Marc Schoenauer, and Michele Sebag. Preference-based policy learning. In *Machine Learning and Knowledge Discovery in Databases*. 2011.
- [2] Riad Akrou, Marc Schoenauer, and Michele Sebag. Interactive robot education. In *European Conference on Machine Learning Workshop*, 2013.
- [3] Ravi Balasubramanian, Ling Xu, Peter D Brook, Joshua R Smith, and Yoky Matsuoka. Physical human interactive guidance: Identifying grasping principles from human-planned grasps. *IEEE Transactions on Robotics*, 2012.
- [4] Jeshua Bratman, Satinder Singh, Jonathan Sorg, and Richard Lewis. Strong mitigation: Nesting search for good policies within search for good reward. In *International Conference on Autonomous Agents and Multiagent Systems*, 2012.
- [5] Maya Cakmak and Andrea L Thomaz. Designing robot learners that ask good questions. In *International conference on Human-Robot Interaction*, 2012.
- [6] Weiwei Cheng, Johannes Fürnkranz, Eyke Hüllermeier, and Sang-Hyeun Park. Preference-based policy iteration: Leveraging preference learning for reinforcement learning. In *Machine Learning and Knowledge Discovery in Databases*. 2011.
- [7] Wei Chu and Zoubin Ghahramani. Preference learning with Gaussian processes. In *International Conference on Machine Learning*, 2005.
- [8] Hao Dang and Peter K Allen. Learning grasp stability. In *International Conference on Robotics and Automation*, 2012.
- [9] Christian Daniel, Gerhard Neumann, and Jan Peters. Learning Sequential Motor Tasks. In *International Conference on Robotics and Automation*, 2013.
- [10] Sam Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In *International Conference on Autonomous Agents and Multiagent Systems*, 2012.
- [11] Marco Dorigo and Marco Colombetti. Robot shaping: Developing autonomous agents through learning. *Artificial intelligence*, 1994.
- [12] Yaakov Engel, Shie Mannor, and Ron Meir. Reinforcement learning with Gaussian processes. In *International Conference on Machine Learning*, 2005.
- [13] Mohammad Ghavamzadeh and Yaakov Engel. Bayesian policy gradient algorithms. *Advances in Neural Information Processing Systems*, 2007.
- [14] Shane Griffith, Kaushik Subramanian, Jonathan Scholz, Charles Isbell, and Andrea L Thomaz. Policy shaping: Integrating human feedback with reinforcement learning. In *Advances in Neural Information Processing Systems*, 2013.
- [15] Harold W Hake and WR Garner. The effect of presenting various numbers of discrete steps on scale reading accuracy. *Journal of Experimental Psychology*, 1951.
- [16] Matthew D Hoffman, Eric Brochu, and Nando de Freitas. Portfolio allocation for Bayesian optimization. In *Conference on Uncertainty in Artificial Intelligence*, 2011.
- [17] Auke Ijspeert, Jun Nakanishi, and Stefan Schaal. Learning Attractor Landscapes for Learning Motor Primitives. In *Advances in Neural Information Processing Systems*. 2003.
- [18] Ashesh Jain, Brian Wojcik, Thorsten Joachims, and Ashutosh Saxena. Learning trajectory preferences for manipulators via iterative improvement. In *Advances in Neural Information Processing Systems*, 2013.
- [19] W Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: The TAMER framework. In *International Conference on Knowledge Capture*, 2009.
- [20] Jens Kober, Betty J. Mohler, and Jan Peters. Learning perceptual coupling for motor primitives. In *Intelligent Robots and Systems*, 2008.
- [21] George Konidaris and Andrew Barto. Autonomous shaping: Knowledge transfer in reinforcement learning. In *International Conference on Machine Learning*, 2006.
- [22] Petar Kormushev, Sylvain Calinon, and Darwin Caldwell. Robot motor skill coordination with EM-based reinforcement learning. In *Intelligent Robots and Systems*, 2010.
- [23] Oliver Kroemer, Renaud Detry, Justus Piater, and Jan Peters. Combining active learning and reactive control for robot grasping. *Robotics and Autonomous Systems*, 2010.
- [24] Manuel Lopes, Francisco Melo, and Luis Montesano. Active learning for reward estimation in inverse reinforcement learning. In *Machine Learning and Knowledge Discovery in Databases*. 2009.
- [25] George A Miller. The magical number seven, plus or minus two: some limits on our capacity for processing information. *Psychological Review*, 1956.
- [26] Andrew Ng and Adam Coates. Autonomous inverted helicopter flight via reinforcement learning. *Experimental Robotics IX*, 1998.
- [27] Andrew Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, 1999.
- [28] Jan Peters, Katharina Mülling, and Yasemin Altun. Relative entropy policy search. In *National Conference on Artificial Intelligence*, 2010.
- [29] Carl E. Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. 2006.
- [30] N. Ratliff, A. Bagnell, and M. Zinkevich. Maximum Margin Planning. In *International Conference on Machine Learning*, 2006.
- [31] Nathan Ratliff, David Silver, and Andrew Bagnell. Learning to Search: Functional Gradient Techniques for Imitation Learning. *Autonomous Robots*, 2009.
- [32] Satinder Singh, Richard L Lewis, Andrew G Barto, and Jonathan Sorg. Intrinsically motivated reinforcement learning: An evolutionary perspective. *Autonomous Mental Development*, 2010.

- [33] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *International Conference on Machine Learning*, 2010.
- [34] Raúl Suárez Feijóo, Jordi Cornellà Medrano, and Máximo Roa Garzón. Grasp quality measures. 2012.
- [35] Andrea L Thomaz and Cynthia Breazeal. Teachable robots: Understanding human teaching behavior to build more effective robot learners. *Artificial Intelligence*, 2008.
- [36] Brian Ziebart, Andrew Maas, Andrew Bagnell, and Anind Dey. Maximum entropy inverse reinforcement learning. In *Conference on Artificial Intelligence*, 2008.