# Skill-based Multi-objective Reinforcement Learning of Industrial Robot Tasks with Planning and Knowledge Integration

Matthias Mayr[1], Faseeh Ahmad[1], Konstantinos Chatzilygeroudis[2], Luigi Nardi[1,3] and Volker Krueger[1]

*Abstract*— In modern industrial settings with small batch sizes it should be easy to set up a robot system for a new task. Strategies exist, e.g. the use of skills, but when it comes to handling forces and torques, these systems often fall short. We introduce an approach that provides a combination of task-level planning with targeted learning of scenario-specific parameters for skill-based systems. We propose the following pipeline: the user provides a task goal in the planning language *PDDL*, then a plan (i.e., a sequence of skills) is generated and the learnable parameters of the skills are automatically identified, and, finally, an operator chooses reward functions and hyperparameters for the learning process. Two aspects of our methodology are critical: (a) learning is tightly integrated with a knowledge framework to support symbolic planning and to provide priors for learning, (b) using multi-objective optimization. This can help to balance key performance indicators (KPIs) such as safety and task performance since they can often affect each other. We adopt a multi-objective Bayesian optimization approach and learn entirely in simulation. We demonstrate the efficacy and versatility of our approach by learning skill parameters for two different contact-rich tasks. We show their successful execution on a real 7-DOF *KUKA-iiwa* manipulator and outperform the manual parameterization by human robot operators.

## I. INTRODUCTION

Industrial environments with expensive and fragile equipment, safety regulations and frequently changing tasks often have special requirements for the behaviour policies that control a robot: First, the trend in industrial manufacturing is to move to smaller batch sizes and higher flexibility of work stations. Reconfiguration needs to be fast, easy and should minimize downtime. Second, it is important to be able to guarantee the performance as well as safety for material and workers. Therefore, it is crucial to be able to understand *what* action is performed *when* and *why*. Finally, in industrial environments digital twins provide a lot of task-relevant information such as material properties and approximate part locations that the robot behavior policies have to consider.

One way to fulfill these criteria is to use systems based on parameterized *skills* [1], [2], [3]. These encapsulated abilities realize semantically defined actions such as moving the robot arm, opening a gripper or localizing an object with vision. State-of-the-art skill-based software architectures can not only utilize knowledge, but also automatically generate
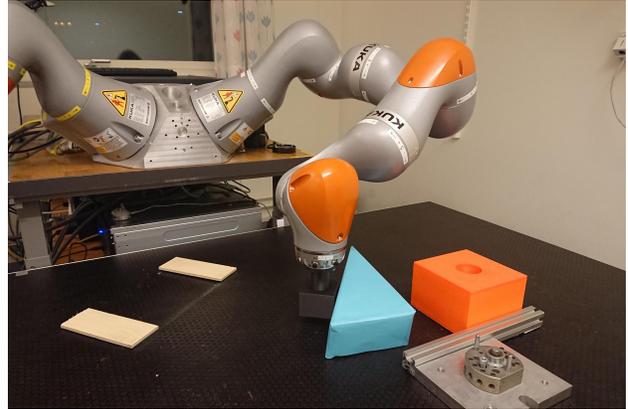


Fig. 1. The robot setup used for the experiments. Wooden boards indicate the start location for the push task. The goal is the corner between the fixture and the box with the hole for the peg task.

plans (skill-sequences) for a given task [4], [5]. The skill-based approach is powerful when knowledge can be modeled and formalized *explicitly* [1], [2]. But it is often limited when it comes to skill parameters of contact-rich tasks that are difficult to reason about. One example are the parameters of a peg insertion search strategy where material properties (e.g. friction) and the robot behavior need to be considered. While it is possible to create a reasoner that follows a set of rules to determine such skill parameters, it is challenging to implement and to maintain.

Another way to handle this is to have operators manually specify and try values for these skill parameters. However, this is a manual process and can be cumbersome.

Finally, it is possible to allow the system to learn by interacting with the environment. However, many policy formulations that allow learning (e.g. artificial networks) have deficiencies which make their application in an industrial domain with the abovementioned requirements challenging. Primarily during the learning phase, dangerous behaviors can be produced and even state-of-the-art RL methods need hundreds of hours of interaction time [6]. Learning in simulation can help to reduce downtime and dangers for the real system. But many policy formulations are black boxes for operators and it can be hard to predict their behavior, which could hinder the trust to the system [7] Moreover, the simulation-to-reality gap [8], [9] is bigger in lower-level control states (i.e. torques), and policies working directly on raw control states struggle to transfer learned behaviors to the real systems [6]. Our policy formulation consisting of behavior trees (BT) with a motion generator [10] has shown to be able to learn interpretable and robust behaviors [11].

[1]Department of Computer Science, Faculty of Engineering (LTH), Lund University, SE 221 00 Lund, Sweden. E-mail: <firstname>.<lastname>@cs.lth.se.

[2]Computer Engineering and Informatics Department (CEID), University of Patras, Greece. E-mail: costashatz@upatras.gr.

[3] Department of Computer Science and Electrical Engineering, Stanford University, CA 94305, USA. E-mail: lnardi@stanford.edu.

The formulation of a learning problem for a given task is often not easy and becomes more challenging if factors such as safety or impact on the workstation environment need to be considered. Multi-objective optimization techniques allow to specify multiple objectives and optimize for them concurrently. This allows operators to select from solutions that are optimal for a certain trade-off between the objectives (usually represented as a set of Pareto-optimal solutions). In order to learn sample-efficient and to support the large variety of skill implementations as well as scenarios, we use gradient-free Bayesian optimization as an optimization method.

In this paper we make the following contributions:

1) We introduce a new method which seamlessly integrates symbolic planning and reinforcement learning for skill-based systems to learn interpretable policies for a given task.

2) A Bayesian multi-objective treatment of the task learning problem, which includes the operator through easy specification of problem constraints and task objectives (KPIs); the set of Pareto-optimal solutions is presented to the operator and their behavior can be inspected in simulation and executed on the real system.

3) We demonstrate our approach on two contact-rich tasks, a pushing task and a peg-in-hole task. We compare it to the outcome of the planner without reasoning, randomly sampled parameter sets from the search space and the manual real-world parameterization process of robot operators. In both tasks our approach delivered solutions that even outperform the ones found by the manual search of human robot operators.

## II. RELATED WORK

### A. Skill-based Systems

Skill-based systems are one way to support a quick setup of a robot system for a new task and to allow re-use of capabilities. There are multiple definitions of the term *skills* in the literature. Some define it as a pure *motion skills* [12] or "hybrid motions or tool operations" [13]. Other work has a broader skill definition [1], [2], [3], [4], [5], [14], [15]. In this formulation, skills can be arbitrary capabilities that change the state of the world and have pre- and postconditions. Their implementation can include motion skills, but also proficiencies such as vision-based localization of objects. In [16] skills are "high-level reusable robot capabilities, with the goal to reduce the complexity and time consumption of robot programming". However, compared to [3] and [14] they do not use pre- and postconditions. In [17], an integrated system for manual creation of *task plans* is presented. It shares the usage of BTs with our approach.

Task planners are used in [1], [2], [4], [5], [14], [18], [16], [13] while [17] lacks such a capability.

In [16] it is suggested that "Machine learning can be performed on the motion level, in terms of adaptation, or can take the form of structured learning on a task/error specification level". However, none of the reviewed work offers a combination task-level planning with learning.

### B. Policy Representation and Learning

An important decision to make when working with manipulators is the type of policy representation and on which level it interfaces with the robot. The latter can strongly influence the learning speed and the quality of the obtained solutions [19], [20]. These choices also influence the form of priors that can be defined and how they are defined [6]. Not many policies combine the aforementioned properties of being a) interpretable, b) paramterizable for the task at hand and c) allow learning or improvement.

The commonly used policy representations for learning systems include radial basis function networks [21], dynamical movement primitives [22], [23] and feed-forward neural networks [21], [24]. In recent years deep artificial neural networks (ANN) seem to become a popular policy. All of them have in common that their final representation can be difficult to interpret. Even if a policy only sets a target pose for the robot to reach, it can be problematic to know how it reacts in all parts of the state space. In contrast to that, [11] suggests to learn interpretable policies based on behavior trees [10] that work explicitly in end-effector space.

### C. Planning and Learning

Symbolic planning is combined with learning in [25], [26], [27], [28]. In [25], the PLANQ-Learning algorithm uses a symbolic planner to shape the reward function based on the conditions defined which are then used by the Q-learner to get an optimal policy with good results on the grid domain. [26] uses the combined symbolic planner with reinforcement learning (RL) in a hierarchical framework to solve complex visual interactive question answering tasks. PEORL [27] integrates symbolic planning and hierarchical reinforcement learning (HRL) to improve performance by achieving rapid policy search and robust symbolic planning in the taxi domain and grid world. SPOTTER [28] uses RL to allow the planning agent to discover the new operators required to complete tasks in Grid World. In contrast to all these approaches, our approach aims towards real-life robotic tasks in an Industry 4.0 setting where a digital twin is available.

In [29], the authors combine symbolic planning with behavior trees (BT) to solve blocks world tasks with a robot manipulator. They use modified Genetic Programming (GP) [30] to learn the structure of the BT. In our approach, we focus on learning the parameters of the skills in the BT and utilize a symbolic planner to obtain the structure of the BT.

## III. APPROACH

Our approach consists of two main components that interact in different stages of the learning pipeline: First, *SkiROS* [14], a skill-based framework for ROS, which represents the implemented skills with BTs, hosts the world model (digital twin), and interacts with the planner. *SkiROS* is also
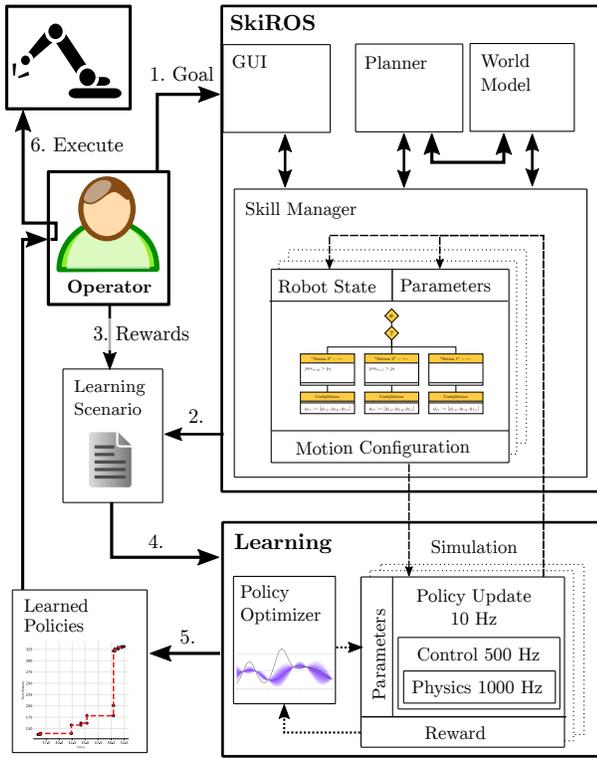
Fig. 2. The architecture of the system that depicts the pipeline: (1) The operator enters the goal state; (2) a learning scenario for the plan is created; (3) rewards and hyperparameters are specified; (4) learning is conducted using the skills and the information in the world model; (5) after policy learning, the operator can choose which policies to execute on the real system (6).



Fig. 3. The BT of the generated plan for the peg insertion task in *eBT format* [32]. Each node has conditions or pre-conditions shown in the upper half and effects or post-conditions shown in the lower half. The *serial start control flow node* ($\rightarrow^*$) executes in a sequence and remembers the successes. The skills have a *parallel-first-success* processor ($<\|FS>$).

the branches according to the implementation of the *control flow nodes* and the return statements of their children. By convention, the signal propagation goes from left to right.

The *sequence* node corresponds to a logical *AND*: it succeeds if all children succeed and fails if one child fails. The *selector*, also called *fallback* node, represents a logical *OR*: If one child succeeds, the remaining ones will not be ticked. It fails only if all children fail. The *parallel* control flow node forwards ticks to all children and fails if one fails. A *decorator* allows to define custom functions. Implementations like *extended Behavior Trees* (eBT) in *SkiROS* [32] add custom processors such as *parallel-first-success* that succeeds if one of the parallel running children succeeds. Leaves of the BT are the *execution nodes* that, when ticked, execute one cycle and output one of the three signals: *success*, *failure* or *running*. In particular, execution nodes subdivide into 1) *action* and 2) *condition* nodes. An action performs its operation iteratively at every tick, returning *running* while it is not done, and *success* or *failure* otherwise. A condition performs an instantaneous operation and returns always *success* or *failure* and never *running*. An example of the BT for the peg insertion task is in Fig. 3.

### B. Planning and Knowledge Integration

The Planning Domain Definition Language (PDDL) [34], [4] is used to formulate the planning problem. We use the *SkiROS* [18] framework that automatically translates a task into a PDDL planning problem by generating domain description and problem instance using the world model. We then use the semantic world model (WM) from *SkiROS* [14] as the knowledge integration framework.

Actions and fluents are obtained by utilizing the predicates that have pre- or post-conditions in the world model. For the problem instance, the objects (robots, arms, grippers, boxes, poses, etc.) in the scene and their initial states (as far as they are known) are used. After getting the necessary

used to execute BTs while learning and to perform tasks on the real system. Second, the learning framework that provides the simulation, the integration with the policy optimizer as well as the reward function definition and calculation. The architecture of the system and the workflow is shown in Figure 2: (1) an operator enters the task goal into a GUI; (2) a plan with the respective learning scenario configuration is generated; (3) an operator complements the scenario with objectives and reward functions; (4) learning is conducted in simulation using the skills and information from the world model; (5) in the multi-objective optimization case, a set of Pareto-optimal solutions is generated and presented to the operator; finally, (6) the operator can select a good solution from this set given the desired trade-off between KPIs and execute it on the real system.

### A. Behavior Trees

A Behavior Tree (BT) [31] is a formalism for plan representation and execution. Like [32], [33], we define it as a directed acyclic graph $G(V, W)$ with $|V|$ nodes and $|W|$ edges. It consists of *control flow nodes* (*processors*), and *execution nodes*. The four basic types of *control flow nodes* are 1) *sequence*, 2) *selector*, 3) *parallel* and 4) *decorator* [33]. A BT always has one initial node with no parents, defined as *Root*, and one or more nodes with no children, called *leaves*. When executing a BT, the *Root* node periodically injects a *tick* signal into the tree. The signal is routed through
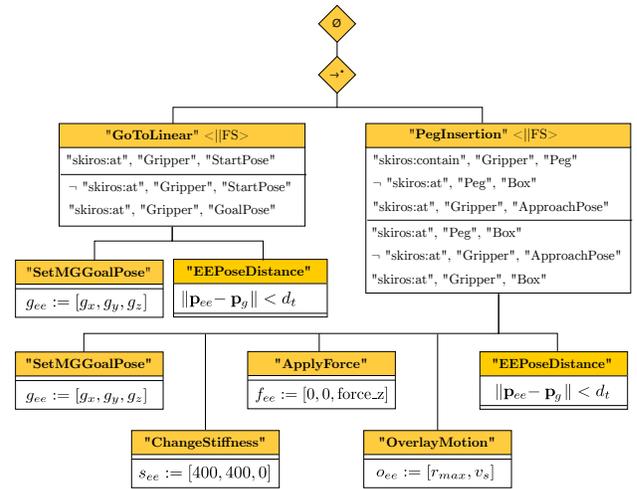
domain description and the problem instance *SkiROS* calls the planner. The goal of the planner is to return a sequence of skills that can achieve the goal conditions of the task. The individual skills are partially parameterized with explicit data from the WM. The WM is aware of the skill parameters that need to be learned for the task at hand and they are automatically identified in the skill sequence.

### C. Policy Optimization

In order to optimize for policy parameters, we adopt the policy search formulation [21], [6], [24]. We formulate a dynamical system in the form:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + M(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\phi}_R), \tag{1}$$

with continuous-valued states $\mathbf{x} \in \mathbb{R}^E$ and actions $\mathbf{u} \in \mathbb{R}^U$. The transition dynamics are modeled by a simulation of the robot and the environment $M(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\phi}_R)$. They are influenced by the domain randomization parameters $\boldsymbol{\phi}_R$.

The goal is to find a policy $\pi, \mathbf{u} = \pi(\mathbf{x}|\boldsymbol{\theta})$ with policy parameters $\boldsymbol{\theta}$ such that we maximize the expected long-term reward when executing the policy for $T$ time steps:

$$J(\boldsymbol{\theta}) = \mathbb{E}\left[\sum_{t=1}^{T} r(\mathbf{x}_t, \mathbf{u}_t)|\boldsymbol{\theta}\right], \tag{2}$$

where $r(\mathbf{x}_t, \mathbf{u}_t)$ is the immediate reward for being in state $\mathbf{x}$ and executing action $\mathbf{u}$ at time step $t$. The discrete switching of branches in the BT and most skills are not differentiable. Therefore, we frame the optimization in Eq. (2) as a black-box optimization and pursue the maximization of the reward function $J(\theta)$ only by using measurements of the function. The optimal reward function to solve the task is generally unknown, and a combination of reward functions is usually used. In the RL literature, this is usually done with a weighted average, that is, $r(\mathbf{x}_t, \mathbf{u}_t) = \sum_i w_i r_i(\mathbf{x}_t, \mathbf{u}_t)$. In this paper, we chose not to use a weighted average of reward functions that represent different objectives (as the optimal combination of weights cannot always be found [35]), but optimize for all objectives concurrently (Sec. III-E) using Bayesian Optimization.

### D. Bayesian Optimization

We consider the problem of finding a global minimizer (or maximizer) of an unknown (black-box) objective function $f$: $\mathbf{s}^* \in \arg\min_{\mathbf{s} \in \mathbb{S}} f(\mathbf{s})$, where $\mathbb{S}$ is some input design space of interest in $D$ dimensions. The problem addressed in this paper is the optimization of a (possibly noisy) function $f : \mathbb{S} \to \mathbb{R}$ with lower and upper bounds on the problem variables. The variables defining $\mathbb{S}$ can be real (continuous), integer, ordinal, and categorical as in [36]. We assume that the function $f$ is in general expensive to evaluate and that the derivatives of $f$ are in general not available. The function $f$ is called black box because we cannot access other information than the output $y$ given an input value $\mathbf{s}$.

This problem can be tackled using Bayesian Optimization (BO) [37]. BO approximates $\mathbf{s}^*$ with a sequence of evaluations, $y_1, y_2, \ldots, y_t$ at $\mathbf{s}_1, \mathbf{s}_2, \ldots, \mathbf{s}_t \in \mathbb{S}$, which maximizes an utility metric, with each new $\mathbf{s}_{t+1}$ depending on the previous function values. BO achieves this by building a probabilistic surrogate model on $f$ based on the set of evaluated points $\{(\mathbf{s}_i, y_i)\}_{i=1}^{t}$. At each iteration, a new point is selected and evaluated based on the surrogate model which is then updated to include the new point $(\mathbf{s}_{t+1}, y_{t+1})$. BO defines an utility metric called the acquisition function, which gives a score to each $\mathbf{s} \in \mathbb{S}$ by balancing the predicted value and the uncertainty of the prediction for $\mathbf{s}$. The maximization of the acquisition function guides the sequential decision making process and the exploration versus exploitation trade-off: the highest score identifies the next point $\mathbf{s}_{t+1}$ to evaluate. BO is a statistically efficient black-box optimization approach when considering the number of necessary function evaluations [38]. It is, thus, especially well-suited to solve problems where we can only perform a limited number of function evaluations, such as the ones found in robotics.

We use the implementation of BO found in *HyperMapper* [36], [39], [40], [41]. Our implementation selects the Expected Improvement (EI) acquisition function [42] and we use uniform random samples as a warm-up strategy before starting the optimization.

### E. Multi-objective Optimization

Let us consider a multiple objectives minimization (or maximization) over $\mathbb{S}$ in $D$ dimensions. We define $f : \mathbb{S} \to \mathbb{R}^p$ as our vector of objective functions $f = (f_1, \ldots, f_p)$, taking $\mathbf{s}$ as input, and evaluating $y = f(\mathbf{s}) + \epsilon$, where $\epsilon$ is a Gaussian noise term. Our goal is to identify the Pareto frontier of $f$, that is, the set $\Gamma \subseteq \mathbb{S}$ of points which are not dominated by any other point, *i.e.,* the maximally desirable $\mathbf{s}$ which cannot be optimized further for any single objective without making a trade-off. Formally, we consider the partial order in $\mathbb{R}^p$: $y \prec y'$ iff $\forall i \in [p], y_i \leqslant y_i'$ and $\exists j, y_j < y_j'$, and define the induced order on $\mathbb{S}$: $\mathbf{s} \prec \mathbf{s}'$ iff $f(\mathbf{s}) \prec f(\mathbf{s}')$. The set of minimal points in this order is the Pareto-optimal set $\Gamma = \{\mathbf{s} \in \mathbb{S} : \nexists \mathbf{s}' \text{ such that } \mathbf{s}' \prec \mathbf{s}\}$. We aim to identify $\Gamma$ with the fewest possible function evaluations using BO. For this purpose we use the *HyperMapper* multi-objective Bayesian optimization which is based on random scalarizations [43].

### F. Motion Generator and Robot Control

The arm motions are controlled in end-effector space by a Cartesian impedance controller. The time varying *reference* or *attractor point* of the end effector $\mathbf{x}_d$ is governed by a motion generator (MG). Given the joint configuration $\mathbf{q}$, we can calculate the end-effector pose $\mathbf{x}_{ee}$ using forward kinematics and obtain an error term $\mathbf{x}_e = \mathbf{x}_{ee} - \mathbf{x}_d$. Together with the joint velocities $\dot{\mathbf{q}}$, the Jacobian $\mathbf{J}(\mathbf{q})$, the configurable stiffness and damping matrices $K_d$ and $D_d$, the task control is formulated as $\tau_c = \mathbf{J}^T(\mathbf{q}) \left(-\mathbf{K}_d \mathbf{x}_e - \mathbf{D}_d \mathbf{J}(\mathbf{q}) \dot{\mathbf{q}}\right)$. Additionally, the task control can be overlaid with commanded generalized forces and torques $\mathbf{F}_{ext} = (f_x\ f_y\ f_z\ \tau_x\ \tau_y\ \tau_z)$: $\tau_{ext} = \mathbf{J}^T(\mathbf{q})\mathbf{F}_{ext}$. We utilize the integration introduced in [10] and used in [11], which proposes to parameterize the MG with movement skills from the BT. The reference pose is shaped by 1) a linear trajectory to a goal point and 2)

overlay motions that can be added to the reference pose as discussed in [10], [11]. E.g. an Archimedes spiral for search.

To make it compliant with the dynamical system in Eq. (1), a new reference configuration of the controller is only generated at every time step $t$. It includes the reference pose, stiffnesses, applied wrench and forms the action $\mathbf{u}$ with a dimension of $U = 19$. The stiffness and applied force are changed gradually at every time step $t$ to ensure a smooth motion. The state space consists of joint positions and joint velocities and is $E = 14$ dimensional. Direct control of the torques of a robot arm requires high update rates and we control the robot arm at 500 Hz based on the current action $\mathbf{u}$, but continuously updated values for $\mathbf{q}$ and $\dot{\mathbf{q}}$. Therefore, from the perspective of Eq. (1), the controller is to be seen as part of the model $M(\mathbf{x}_t, \mathbf{u}_t)$.

We assume a human-robot collaborative workspace with fragile objects. Therefore, the stiffnesses and applied forces are to be kept to a minimum and less accuracy than e.g. high-gain position-controlled solutions is to be expected.

## IV. Experiments

In our experiments we use a set of pre-defined skills that are part of a skill library. In order to solve a task, the planner determines a sequence that can achieve the goal condition of the task. This skill sequence is also automatically parameterized to the extend possible, e.g. the goal pose of a movement. We evaluate our system in two contact-rich scenarios that are shown in Fig. 1: A) pushing an object with uneven weight distribution to a goal pose and B) inserting a peg in a hole with a $1.5\,\text{mm}$ larger radius. Pure planning-based solutions for both these tasks have a poor performance in reality (Fig. 5).

As a baseline we invited six robot operators to manually parameterize the skills for the tasks. Their main objective is to find a parameter set that robustly solves the task. As an additional objective they were asked to minimize the impact of the robot arm and its tool on the environment as long as it does not affect the first objective.

The robot arm used for the physical evaluation is a 7-degree-of-freedom (DOF) *KUKA iiwa* arm controlled by a Cartesian impedance controller (Sec. III-F).

### A. Reward Functions

For each task, we utilize a set of reward functions parameterized for the learning scenario configuration. Each configured reward has an assigned objective and can be weighted against other rewards. Each experiment uses a subset of the following reward functions:

*1) Task completion:* A fixed reward is assigned when the BT returns success upon task completion.

*2) End-effector distance to a box:* We use a localized reward to attract the end effector towards the goal location $r_h(\mathbf{x}) = \left(2\left(d(\mathbf{p}_{ee,\mathbf{x}}, \mathbf{p}_h) + d_o\right)\right)^{-1}$, where $d_o$ is the distance offset and $d(\mathbf{p}_{ee,\mathbf{x}}, \mathbf{p}_h)$ is the shortest distance function between the end effector and the box.

*3) Applied wrench:* This reward calculates the cumulative forces applied by the end effector on the environment.

Reward functions 4-6 share a common operation of computing an exponential function of the calculated metric to obtain the reward as used in ([44], [24]) $r(d_m) = \exp\left(-\frac{1}{2\sigma_w^2}(d_m + d_o)\right)$, where $\sigma_w$ is a configurable width, $d_o$ is a distance offset and $d_m$ is the input metric.

*4) End-effector distance to a goal:* This reward uses distance between the end effectors current pose and goal pose to calculate the input metric $d_{ee,g} = \|\mathbf{p}_{ee,\mathbf{x}} - \mathbf{p}_g\|$

*5) End-effector-reference-position distance:* This reward uses the distance between the end effectors reference pose (Sec. III-F) and its current pose to calculate the input metric $d_{ee,d} = \|\mathbf{p}_{ee,\mathbf{x}} - \mathbf{x_d}\|$

*6) Object-pose divergence:* This reward uses the translational and angular distance between the object's goal pose and its current pose.

### B. Push Task

The push task starts by specifying the goal in the *SkiROS* Graphical User Interface (GUI) as: (skiros:at skiros:ObjectToBePushed-1 skiros:ObjectGoalPose-1). *SkiROS* calls the planner to generate a plan given all the available skills. The plan consists of two skills: 1) $GoToLinear$ skill and 2) $Push$ skill. The first skill moves the end effector from its current location to the *approach pose* of the object. This *approach pose* is defined in the WM and needs to be reached before interacting with the object.

The push skill then moves the end effector to the object's geometric centre with an optional offset in the horizontal ($x$) and ($y$) directions. Once the end effector reaches it, the motion generator executes a straight line to the (modified) target location.

The push task is formulated as a multi-objective task. It also has two objectives, 1) success and 2) applied force. The first objective has three associated rewards: 1) object position difference from goal position, 2) object orientation difference from goal orientation, and 3) end-effector distance to the goal location. The second objective accumulates the Cartesian distance between the end-effector reference pose and the actual end-effector pose as a measure of the force applied by the controller. The learnable parameters in this task are offsets in the horizontal ($x$) and ($y$) direction of both the push skill's start and goal locations. An offset of the start location allows the robot to push from a particular point from the side of the object. Together with the offsets on the goal position, these learnable parameters collectively define the trajectory of the push.

The object to be pushed has a height of $0.07\,\text{m}$ and is an orthogonal triangle in the horizontal dimensions ($x$) and ($y$). It has a length of $0.15\,\text{m}$ and $0.3\,\text{m}$ and it weights $2.5\,\text{kg}$. For this task we use a square-shaped peg for pushing with a side length of $0.07\,\text{m}$ and a height of $0.05\,\text{m}$. Start and goal locations are $\approx 0.43\,\text{m}$ apart and are rotated by $26\,deg$. We define success if the translational and rotational difference of the object w.r.t the goal is less than $0.01\,\text{m}$ and
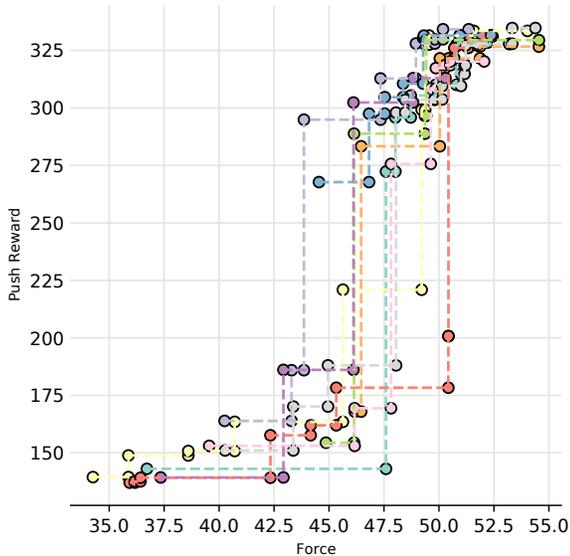
Fig. 4. Pareto front of the push task. Each experiment has a different color and each point represents a Pareto-optimal solution. It shows that higher rewards for pushing require higher interaction forces with the environment.

5 $deg$, respectively. We learn for 400 iterations and repeat the experiment 10 times. In order to obtain solutions that are robust enough to translate to the real system, we apply domain randomization. Each parameter set is evaluated in 7 worlds. Each execution uniformly samples one out of the four start positions for the robot arm. Furthermore, we vary the location of the object and the goal in the horizontal $(x)$ and $(y)$ directions by sampling from a Gaussian distribution with a standard deviation of 7 mm.

We compare the learned solutions with (a) the outcome of a direct planner solution without any offset on the start and goal pose while pushing, (b) ten sets of random parameters from the search space and (c) the policies that are parameterized by the robot operators. We evaluated on the four start configurations used for learning as well as on two additional unknown ones. The results are shown in Fig. 5a.

The results of a multi-objective optimization are parameters found along a Pareto front (Sec. III-E, see Fig. 4). It contained 8.3 points on average, of which some minimize the impact on the environment to an extent that the push is not successful. An operator can choose a solution that is a good compromise between the success of the task on the real system and the force applied on the environment. The performance of one of the solutions that existed on the Pareto front is shown in Fig. 5.

Furthermore, we asked six robot operators to find values for the learnable parameters of the skill sequences. They were given the same start positions used for learning and were given a script to reset the arm to a start position of their choice. They could experiment with the system until they decided that their parameter set fulfills the criteria. Their final parameter set that was also evaluated on the known and unknown start configurations. On average the operators spent $(16.3 \pm 6.4)$ min and executed $11.1 \pm 3.0$ trials on the system to configure this task. Four out of the six

operators found solutions that achieved the task from every start state. However, two of the operators' final parameters only achieved success rates of $50\%$ and $16.66\%$.

### C. Peg-in-Hole Task

The PDDL goal of the peg insertion task is (`skiros:at skiros:Peg-1 skiros:BoxWithHole-1`). The BT that is generated by the planner is shown in Fig. 3 and uses two skills: 1) $GoToLinear$ skill and 2) $PegInsertion$ skill. The first skill moves the end effector from its current location to the *approach pose* of the hole. Once it is reached, the peg insertion procedure starts.

The $PegInsertion$ skill starts when the end effector reaches the approach pose of the box. It uses four separate *SkiROS* primitive skills to 1) set the stiffness of the end effector to zero in $(z)$ direction, 2) apply a downward force in $(z)$ direction, 3) configure the center of the box as a goal and 4) additionally apply an overlaying circular search motion on top of the reference pose of the end effector as described in [11]. The BT returns success only if the peg is inserted into the box hole by more than $0.01$ m.

We model the peg insertion as a multi-objective and multi-reward task. There are two objectives of the task, 1) successful insertion and 2) applied force. To assess the efficacy of the first objective, we use three rewards, 1) success of the BT, 2) peg distance to the hole, and 3) peg distance to the box. For the second objective, we use a single reward that measures the total force applied by the peg. There are three learnable parameters in this task, 1) downward force applied by the robot arm, 2) radius of the overlay search motion and 3) path velocity of the overlay search motion.

We learn for 400 iterations in the simulation and repeat this experiment 10 times. To increase the robustness of the solutions we use domain randomization and evaluate each parameter configuration in 7 worlds. We vary the location of the box by sampling from a Gaussian distribution with a standard deviation of 7 mm and uniformly sample one out of 5 start configurations of the robot arm. We compare the performance of the learned policies with (1) the outcome of the planner without a parameterized search motion, (2) randomly chosen parameter configurations from the parameter search space used for learning and (3) policies that are parameterized by human operators (see Fig. 5b).

The learned Pareto-optimal configurations consist of 6.1 points on average. We evaluated the insertion success using the 5 known and additional 10 unknown start configurations of the robot (Fig. 5b).

To find policies for this task, the human operators took $(31.8 \pm 10.9)$ min and executed $39 \pm 14$ trials on the system. However, compared to the randomly sampled policies the average insertion rate only increased from $41\%$ to $52.2\%$. This is much lower than the average insertion rate of $96\%$ of the best learned policies as shown in box four, Fig. 5b. Furthermore, the average force that was chosen by the operators compared to the learned policies was $16.6\%$ higher. Finally, the successful insertions by the learned policies were also $18.1\%$ faster. Therefore, the learned policies outperformed
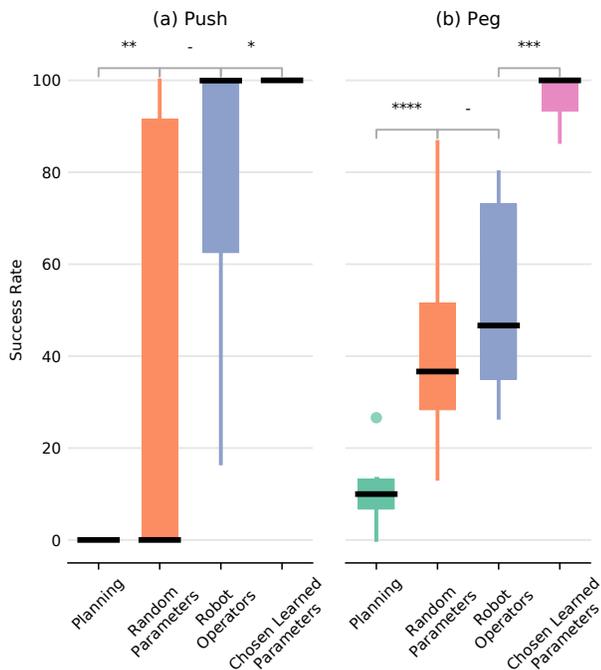
Fig. 5. The success rates of both experiments. The box plots show the median (black line) and interquartile range ($25^{th}$ and $75^{th}$ percentile); the lines extend to the most extreme data points not considered outliers, and outliers are plotted individually. The number of stars indicates that the p-value of the Mann-Whitney U test is less than 0.1, 0.05, 0.01 and 0.001 respectively.

the human operators in both objectives while also producing more reliable results.

## V. CONCLUSION

In this paper we proposed a method for effectively combining task-level planning with learning to solve industrial contact-rich tasks. Our method leverages prior information and planning to acquire *explicit* knowledge about the task, whereas it utilizes learning to capture the *tacit* knowledge, i.e., the knowledge that is hard to formalize and which can only be captured through actual interaction. We utilize behavior trees as an interpretable policy representation that is suitable for learning and leverage domain randomization for learning in simulation. Finally, we formulate a multi-objective optimization scheme so that (1) we handle conflicting rewards adequately, and (2) an operator can choose a policy from the Pareto front and thus actively participate in the learning process.

We evaluated our method on two scenarios using a real *KUKA* 7-DOF manipulator: (a) a pushing task, and (b) a peg insertion task. Both tasks are contact-rich and naïve planning fails to solve them. The approach was able to outperform the baselines including the manual parameterization by robot operators.

For future work we are looking into multi-fidelity learning that can leverage a small amount of executions on the real system to complement the learning in simulation. Furthermore, the use of parameter priors for the optimum seems a promising direction to guide the policy search and make it more efficient.
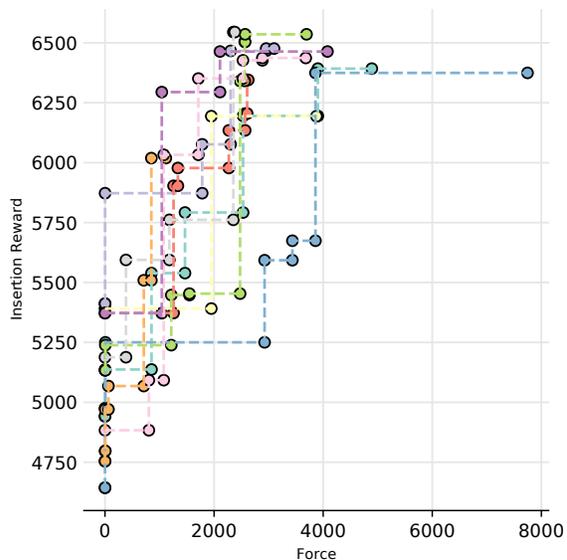


Fig. 6. Pareto front of the peg task. Each experiment has a different color. The goal is to maximize insertion reward while minimizing the interaction forces.

## APPENDIX

The implementation and the supplemental video are available at:

`https://sites.google.com/ulund.org/SkiREIL`

## REFERENCES

[1] V. Krueger, F. Rovida, B. Grossmann, R. Petrick, M. Crosby, A. Charzoule, G. Martin Garcia, S. Behnke, C. Toscano, and G. Veiga, "Testing the vertical and cyber-physical integration of cognitive robots in manufacturing," *Robotics and Computer-Integrated Manufacturing*, vol. 57, pp. 213–229, 2019.

[2] V. Krueger, A. Chazoule, M. Crosby, A. Lasnier, M. R. Pedersen, F. Rovida, L. Nalpantidis, R. Petrick, C. Toscano, and G. Veiga, "A vertical and cyber–physical integration of cognitive robots in manufacturing," *Proceedings of the IEEE*, vol. 104, no. 5, pp. 1114–1127, 2016.

[3] S. Bøgh, O. S. Nielsen, M. R. Pedersen, V. Krüger, and O. Madsen, "Does your robot have skills?" in *Proceedings of the 43rd international symposium on robotics*. VDE Verlag GMBH, 2012.

[4] M. Crosby, F. Rovida, V. Krueger, and R. P. A. Petrick, "Integrating mission and task planning in an industrial robotics framework," in *Proceedings of the Twenty-Seventh International Conference on Automated Planning and Scheduling (ICAPS 2017)*. AAAI, 2017.

[5] F. Rovida, V. Krüger, C. Toscano, G. Veiga, M. Crosby, and R. Petrick, "Planning for sustainable and reliable robotic part handling in manufacturing automation," in *Proceedings of the Twenty-Sixth International Conference on Automated Planning and Scheduling (ICAPS 2016)*, 2016.

[6] K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, and J.-B. Mouret, "A survey on policy search algorithms for learning robot controllers in a handful of trials," *IEEE Transactions on Robotics*, vol. 36, no. 2, pp. 328–347, 2019.

[7] M. Edmonds, F. Gao, H. Liu, X. Xie, S. Qi, B. Rothrock, Y. Zhu, Y. N. Wu, H. Lu, and S.-C. Zhu, "A tale of two explanations: Enhancing human trust by explaining robot behavior," *Science Robotics*, vol. 4, no. 37, p. eaay4663, 2019.

[8] S. Koos, J.-B. Mouret, and S. Doncieux, "The transferability approach: Crossing the reality gap in evolutionary robotics," *IEEE Transactions on Evolutionary Computation*, vol. 17, no. 1, pp. 122–145, 2012.

[9] J.-B. Mouret and K. Chatzilygeroudis, "20 years of reality gap: a few thoughts about simulators in evolutionary robotics," in *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 2017, pp. 1121–1124.

[10] F. Rovida, D. Wuthier, B. Grossmann, M. Fumagalli, and V. Krüger, "Motion Generators Combined with Behavior Trees: A Novel Approach to Skill Modelling," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 5964–5971.

[11] M. Mayr, K. Chatzilygeroudis, F. Ahmad, L. Nardi, and V. Krueger, "Learning of Parameters in Behavior Trees for Movement Skills," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021.

[12] T. Hasegawa, T. Suehiro, and K. Takase, "A model-based manipulation system with skill-based execution in unstructured environment," in *Fifth International Conference on Advanced Robotics' Robots in Unstructured Environments*. IEEE, 1991, pp. 970–975.

[13] U. Thomas, B. Finkemeyer, T. Kroger, and F. M. Wahl, "Error-tolerant execution of complex robot tasks based on skill primitives," in *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*, vol. 3. IEEE, 2003, pp. 3069–3075.

[14] F. Rovida, M. Crosby, D. Holz, A. S. Polydoros, B. Großmann, R. P. Petrick, and V. Krüger, "SkiROS-A skill-based robot control platform on top of ROS," in *Studies in Computational Intelligence*, 2017, vol. 707, pp. 121–160.

[15] U. Thomas, G. Hirzinger, B. Rumpe, C. Schulze, and A. Wortmann, "A new skill based robot programming language using uml/p statecharts," in *2013 IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 461–466.

[16] M. Stenmark, J. Malec, K. Nilsson, and A. Robertsson, "On distributed knowledge bases for robotized small-batch assembly," *IEEE Transactions on Automation Science and Engineering*, vol. 12, no. 2, pp. 519–528, 2015.

[17] C. Paxton, A. Hundt, F. Jonathan, K. Guerin, and G. D. Hager, "Costar: Instructing collaborative robots with behavior trees and vision," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 564–571.

[18] F. Rovida, B. Grossmann, and V. Krüger, "Extended behavior trees for quick definition of flexible robotic tasks," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 6793–6800.

[19] P. Varin, L. Grossman, and S. Kuindersma, "A comparison of action spaces for learning manipulation tasks," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 6015–6021.

[20] R. Martín-Martín, M. A. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg, "Variable impedance control in end-effector space: An action space for reinforcement learning in contact-rich tasks," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 1010–1017.

[21] M. P. Deisenroth, G. Neumann, and J. Peters, "A Survey on Policy Search for Robotics," *Foundations and Trends® in Robotics*, vol. 2, no. 1–2, pp. 1–142, 2013. [Online]. Available: https://www.nowpublishers.com/article/Details/ROB-021

[22] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal, "Dynamical movement primitives: learning attractor models for motor behaviors," *Neural computation*, vol. 25, no. 2, pp. 328–373, 2013.

[23] A. Ude, A. Gams, T. Asfour, and J. Morimoto, "Task-specific generalization of discrete and periodic dynamic movement primitives," *IEEE Transactions on Robotics*, vol. 26, no. 5, pp. 800–815, 2010.

[24] K. Chatzilygeroudis, R. Rama, R. Kaushik, D. Goepp, V. Vassiliades, and J.-B. Mouret, "Black-box data-efficient policy search for robotics," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 51–58.

[25] M. Grounds and D. Kudenko, "Combining reinforcement learning with symbolic planning," in *Adaptive Agents and Multi-Agent Systems III. Adaptation and Multi-Agent Learning*. Springer, 2005, pp. 75–86.

[26] D. Gordon, D. Fox, and A. Farhadi, "What should i do now? mar-

[27] rying reinforcement learning and symbolic planning," *arXiv preprint arXiv:1901.01492*, 2019.

[27] F. Yang, D. Lyu, B. Liu, and S. Gustafson, "Peorl: Integrating symbolic planning and hierarchical reinforcement learning for robust decision-making," *arXiv preprint arXiv:1804.07779*, 2018.

[28] V. Sarathy, D. Kasenberg, S. Goel, J. Sinapov, and M. Scheutz, "Spotter: Extending symbolic planning operators through targeted reinforcement learning," in *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 2021.

[29] J. Styrud, M. Iovino, M. Norrlöf, M. Björkman, and C. Smith, "Combining planning and learning of behavior trees for robotic assembly," *arXiv preprint arXiv:2103.09036*, 2021.

[30] J. R. Koza and J. R. Koza, *Genetic programming: on the programming of computers by means of natural selection*. MIT press, 1992, vol. 1.

[31] M. Colledanchise and P. Ögren, *Behavior Trees in Robotics and AI: An Introduction*. Chapman & Hall/CRC Press, 2017.

[32] F. Rovida, B. Grossmann, and V. Krüger, "Extended behavior trees for quick definition of flexible robotic tasks," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 6793–6800.

[33] A. Marzinotto, M. Colledanchise, C. Smith, and P. Ögren, "Towards a unified behavior trees framework for robot control," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 5420–5427.

[34] M. Fox and D. Long, "Pddl2. 1: An extension to pddl for expressing temporal planning domains," *Journal of artificial intelligence research*, vol. 20, pp. 61–124, 2003.

[35] R. Kaushik, K. Chatzilygeroudis, and J.-B. Mouret, "Multi-objective model-based policy search for data-efficient learning with sparse rewards," in *Conference on Robot Learning*. PMLR, 2018, pp. 839–855.

[36] L. Nardi, D. Koeplinger, and K. Olukotun, "Practical design space exploration," in *International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems*, 2019.

[37] P. I. Frazier, "A tutorial on bayesian optimization," 2018.

[38] E. Brochu, V. M. Cora, and N. De Freitas, "A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning," *arXiv preprint arXiv:1012.2599*, 2010.

[39] L. Nardi, B. Bodin, S. Saeedi, E. Vespa, A. J. Davison, and P. H. Kelly, "Algorithmic performance-accuracy trade-off in 3d vision applications using hypermapper," in *International Parallel and Distributed Processing Symposium Workshops*, 2017.

[40] B. Bodin, L. Nardi, M. Z. Zia, H. Wagstaff, G. Sreekar Shenoy, M. Emani, J. Mawer, C. Kotselidis, A. Nisbet, M. Lujan *et al.*, "Integrating algorithmic parameters into benchmarking and design space exploration in 3d scene understanding," in *International Conference on Parallel Architectures and Compilation*, 2016.

[41] A. Souza, L. Nardi, L. B. Oliveira, K. Olukotun, M. Lindauer, and F. Hutter, "Bayesian optimization with a prior for the optimum," in *Machine Learning and Knowledge Discovery in Databases. Research Track*, N. Oliver, F. Pérez-Cruz, S. Kramer, J. Read, and J. A. Lozano, Eds. Cham: Springer International Publishing, 2021, pp. 265–296.

[42] J. Mockus, V. Tiesis, and A. Zilinskas, "The application of bayesian methods for seeking the extremum," *Towards global optimization*, vol. 2, no. 117-129, p. 2, 1978.

[43] B. Paria, K. Kandasamy, and B. Póczos, "A flexible framework for multi-objective bayesian optimization using random scalarizations," in *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI*, A. Globerson and R. Silva, Eds., 2019, p. 267.

[44] M. P. Deisenroth and C. E. Rasmussen, "PILCO: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ser. ICML'11. Omnipress, 2011, pp. 465–472.