

# Q-CP: Learning Action Values for Cooperative Planning

Francesco Riccio<sup>1</sup>

Roberto Capobianco<sup>1</sup>

Daniele Nardi<sup>1</sup>

**Abstract**—Research on multi-robot systems has demonstrated promising results in manifold applications and domains. Still, efficiently learning an effective robot behaviors is very difficult, due to unstructured scenarios, high uncertainties, and large state dimensionality (e.g. hyper-redundant and groups of robot). To alleviate this problem, we present *Q-CP* a cooperative model-based reinforcement learning algorithm, which exploits action values to both (1) guide the exploration of the state space and (2) generate effective policies. Specifically, we exploit *Q*-learning to attack the curse-of-dimensionality in the iterations of a Monte-Carlo Tree Search. We implement and evaluate *Q-CP* on different stochastic cooperative (general-sum) games: (1) a simple cooperative navigation problem among 3 robots, (2) a cooperation scenario between a pair of KUKA YouBots performing hand-overs, and (3) a coordination task between two mobile robots entering a door. The obtained results show the effectiveness of *Q-CP* in the chosen applications, where action values drive the exploration and reduce the computational demand of the planning process while achieving good performance.

## I. INTRODUCTION

Robots have to show robust behaviors to complete cognitive tasks in different scenarios, such as service robotics and uncontrolled industrial environments. However, deciding the best action to perform, is a complex task due to unpredictabilities of the physical world, uncertainties in the observations, continuous state spaces and rapid explosions of the state dimensionality. This is especially true in multi-robot scenarios such as coordinated navigation and cooperative manipulation, where each robot has to represent the state of the environment, and estimate other robots' states in order to determine the most effective action to perform and, to achieve an individual or collective goal. In these scenarios, the number of robots induces a large state-space, where generalization and policy generation become particularly difficult tasks to tackle. The generalization problem is typically addressed with function approximators (e.g. neural networks), but they do not allow for the use of prior knowledge, which can be inefficient and lead to dangerous situations. To overcome this issue, decision theoretical planning techniques, such as Monte-Carlo tree search [1], have been used to embed prior knowledge in learning problems. Nevertheless, they show difficulties in relating similar states (i.e. nodes of the search tree) [2]. Here, we focus on the problem of cooperative general-sum stochastic games [3], where each agent runs its own learning process. We attack the generalization problem in policy generation by enhancing the Upper Confidence Tree (UCT) algorithm [4] – a variant

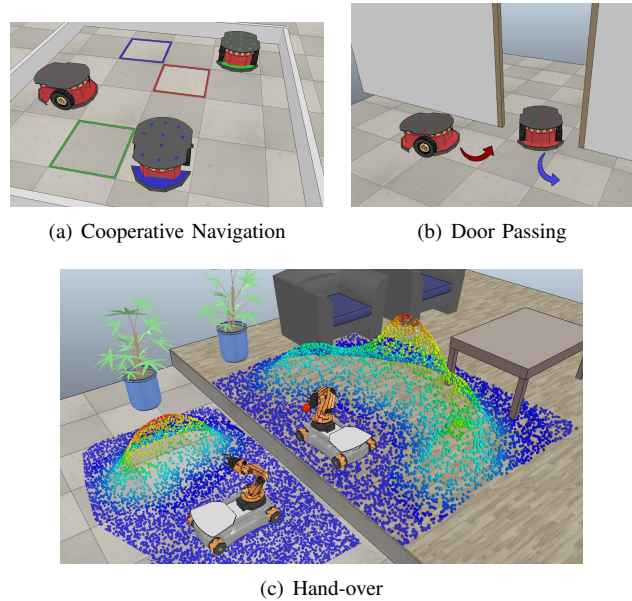


Fig. 1. Multi-robot behaviors learned with *Q-CP*. In each of the applications, the robots exploit generated policies to succeed in: (1) a *cooperative navigation* task (1(a)), where each robot has to reach the square with the matching color; (2) *hand-over* task (1(c)) where a robot has to pass a ball among each other. The two Gaussian functions model the current spatial distribution of two action values generated during learning, i.e. *move-left* for the robot above and *move-right* for the one below; and (3) a *door passing* task, where two robots have to coordinate in order to traverse a narrow passage in opposite directions.

of Monte-Carlo Tree Search – with an external action-value function approximator, that selects admissible actions and consequently drives the node-expansion phase during episode simulation.

In this paper, we introduce *Q-CP*, Q-value based Cooperative Planning, to iteratively learn cooperative policies in stochastic games characterized by discrete action spaces and large state spaces. To model action values, we use regression on a mixture of Gaussians [5], that is iteratively refined by aggregating new training samples [6], after each search completion (i.e., at every  $t$ ). Specifically, *Q-CP* generates robot action policies by running multiple Monte Carlo tree searches [7] and incrementally collecting new samples that are used to improve *Q*-value (i.e., action value) estimates. These are then used to select admissible actions during the search process itself. In our experiments, we present a set of fully collaborative games – where all the robots have identical reward functions – and we therefore choose to not model joint actions [8], [9], thus avoiding further explosions of the UCT search tree.

<sup>1</sup>All authors are with the Department of Computer, Control, and Management Engineering “Antonio Ruberti”, Sapienza University of Rome, via Ariosto 25, Rome, 00185, Italy [lastname@dis.uniroma1.it](mailto:lastname@dis.uniroma1.it)

We aim at demonstrating that  $Q$ -CP can efficiently be used to generalize the policy and restrict the search space to support learning both at an individual and collective level and thus, to enable teamwork among robots that operate with a common goal. To this end, we specifically address applications, where cooperation among robots is necessary to achieve a goal and to improve task performance, as illustrated in Fig. 1. The experimental evaluation shows the effectiveness of explicitly embedding action values in the search process, resulting in a better generalization and a reduction in the computational demand of the learning process. Our main contributions consist in (1) a novel integration of Monte-Carlo tree search, data aggregation and  $Q$ -learning, that enables good performance with large search spaces, (2) an extension of TD-search [2] that constructs upper confidence bounds on the value function to select actions optimistically, and (3) a reduction of the curse-of-dimensionality that is obtained by means of such a focused exploration. These improvements make our approach more practical and suitable in difficult robotics applications, where the lack of training examples is often a limiting condition.

The remainder of the paper is organized as follows. Section II reports recent research on reinforcement learning and multi-robot planning, and Section III describes the  $Q$ -CP algorithm. Finally, Section IV describes our experimental setup, as well as the obtained results, and Section V closes the paper with final remarks, limitations and future directions.

## II. RELATED WORK

Especially in robotics applications where robots are expected to show robust and adaptable behaviors, generating an effective policy is often complex and resource intensive, due to large state spaces, unstructured environments and unpredictabilities. In particular, techniques that are traditionally adopted in planning or learning are computationally demanding and time consuming [1], [10]. To tackle these problems, several approaches initialize robot policies with reasonable behaviors (e.g., through the aid of expert demonstrations [11], [12]), or exploit hierarchical representations like *options* [13], [14] and MAX-Q decompositions [15]. Unfortunately, the applicability of these methods in the multi-agent domain is rather limited and unexplored [16]. Conversely, in this setting, several approaches use Q-learning and decentralized POMDPs to solve a multitude of problems, such as cooperative navigation in 2D grid worlds [17], [18], the *two agent-tiger* and a *box pushing* problems [19], and *prey vs. predator* games in a grid environment [20]. These methods generally require a considerable number of training examples, and are not easily applicable to robotics domains. Besides being affected by large state spaces, traditional decision theoretic planning methods, such as Monte-Carlo tree search, do not provide generalization capabilities, that are required in common robotics problems, where large portions of the state space are rarely or never explored. In robotics, however, generalization may result to be inefficient, slow and can lead to critical situations due to its limitation

in embedding prior knowledge, which is typically available to Monte-Carlo techniques [1].

We address generalization at learning time by introducing  $Q$ -CP, an iterative algorithm for policy generation.  $Q$ -CP allows a team of agents to learn action values, that are used for focused exploration. To this end, we rely on previous literature [21], [22] and we approximate the  $Q$  function using probability densities represented as a mixture of Gaussians [5]. In fact, this type of approximator provides the flexibility, expressiveness and generality of non-parametric function approximators, while retaining the additional benefit of quantifying uncertainty.

Similar to TD-search [2], we aim at reducing the variance of value estimates during the search procedure by means of temporal difference learning. However, we improve our model by preserving the selective search of Monte-Carlo algorithms, when bootstrapping is ongoing. Specifically,  $Q$ -CP extends TD-search by constructing upper confidence bounds on the value function, and by selecting optimistically with respect to those, instead of performing  $\epsilon$ -greedy exploration. In this way,  $Q$ -CP explores more informative portions of the search space [23], [24], and supports the generation of competitive policies – with the additional benefit of a reduction in (1) number of simulations (or expanded nodes), and (2) exploration of the search space – that alleviates the curse-of-dimensionality.

Summarizing, our approach enables a team of robots to generate effective policies directly from simulation, with a reduced training set and number of iterations. Such a feature makes  $Q$ -CP also appealing in robotics, where existing approaches to multi-agent learning cannot be extended nor easily applied.

## III. $Q$ -CP

In this section we present  $Q$ -CP under the stochastic game framework – an extension of Markov Decision Processes to the multi-agent scenario. A stochastic game is a tuple  $(n, \mathcal{S}, \mathbf{A}_{1:n}, \mathcal{T}, R_{1:n})$ , where  $n$  is the number of agents,  $\mathcal{S}$  is the set of states of the environment,  $\mathbf{A}_j$  represents the set of discrete actions of agent  $j$ ,  $\mathcal{T} : \mathcal{S} \times \mathbf{A} \times \mathcal{S} \rightarrow [0, 1]$  is the stochastic transition function that models the probabilities of transitioning from state  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$  when an action is taken from the joint action space  $\mathbf{A} : \mathbf{A}_1 \times \dots \times \mathbf{A}_n$ , and  $R_j : \mathcal{S} \times \mathbf{A} \rightarrow \mathbb{R}$  is the reward function of agent  $i$ . In this setting, decisions are represented through agent policies  $\pi_j$ , that define the behavior of each agent  $j$  by mapping states to actions. Given a stochastic game, the goal of each agent consists in finding a policy  $\pi_j(s)$  that maximizes its future reward with a discount factor  $\gamma$ .

Our goal is to enable teamwork among robots that operate with a common goal. For this reason, we restrict to fully collaborative games [8], where all the agents have identical reward functions.  $Q$ -CP iteratively evolves by (1) running, for each agent an UCT search, where admissible actions are selected through  $Q$ -value estimates and (2) incrementally collecting new samples that are used to improve  $Q$ -value estimates. In this section, we first describe how we approximate

the  $Q$  function; then, we present a detailed explanation of our algorithm.

### A. Preliminaries

We choose to approximate the  $Q$  function using probability densities in the form of a mixture of  $K$  Gaussians (i.e., Gaussian Mixture Models – GMMs), with  $K$  determined in an adaptive manner. The use of such a non-parametric function approximator allows, in fact, to obtain flexible, expressive and general representations, while retaining the additional benefit of quantifying uncertainty. To this end, in fact, we integrate the approach in [5] with a data aggregation [6] procedure, where a dataset of samples is iteratively collected and aggregated. Specifically, at each iteration  $i$  of  $Q$ -CP we collect a dataset  $D_j^i = \{x\}$  of sample state-action pairs and  $Q$ -values  $x = (s, a_j, Q_j^i(s, a_j))$  experienced by agent  $j$ , together with the corresponding rewards  $r$  and transitioned states  $s'$  resulting from the executed joint action. The values  $Q_j^i$  are determined, at each iteration  $i$ , according to the  $Q$ -learning update rule

$$Q_j^i(s, a_j) = \hat{Q}_j^{i-1}(s, a_j) + \alpha (r + \gamma \max_{a'_j} \hat{Q}_j^{i-1}(s', a'_j) - \hat{Q}_j^{i-1}(s, a_j)), \quad (1)$$

where  $\alpha$  is the learning rate,  $\hat{Q}_j^{i-1}$  is the function approximation, and  $Q_j^0(s, a_j) = 0$ . As we discuss later, the function approximation  $\hat{Q}$  is learned over an aggregated dataset  $D_j^{0:i} = \{\cup D_j^d | d = 0 \dots i\}$ . More specifically, the aggregated dataset is used to estimate a probability density function (pdf) in the space of states, actions and  $Q$ -values.

The pdf depends on the set of parameters  $\Theta = \{\pi_1, \mu_1, \Sigma_1, \dots, \pi_K, \mu_K, \Sigma_K\}$ , where  $\pi_k$  is the prior,  $\mu_k$  the mean and  $\Sigma_k$  the covariance matrix of a Gaussian of dimensionality  $G$  in the canonical form

$$N(x, \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^G |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}. \quad (2)$$

The parameters  $\Theta$  are estimated over the dataset  $D_j^{0:i}$ , and they are obtained as the result of a standard Expectation-Maximization [25] procedure, initialized using the k-means algorithm [26] to avoid bad local optima. The number of components  $K$  of the GMM is selected to minimize the Bayesian Information Criterion [27] over a testing portion of the dataset.

Using  $\Theta$ , the pdf of a sample can be computed as

$$p(x; \Theta) = \sum_{k=1}^K \pi_k N(x, \mu_k, \Sigma_k), \quad (3)$$

while the value  $\hat{Q}(s, a) \approx Q(s, a)$  is obtained as  $\hat{Q}(s, a) = \mathbb{E}[Q|s, a] = \mu(Q|s, a)$ . The approximated  $Q$ -value is the result of a Gaussian Mixture Regression

$$\begin{aligned} \hat{Q}(s, a) &= \mu(Q|s, a) = \sum_{k=1}^K \beta(s, a)_k \mu_k(Q|s, a), \\ \sigma^2(Q|s, a) &= \sum_{g=1}^G \beta(s, a)_g y_g(s, a) - \mu^2(Q|s, a), \end{aligned} \quad (4)$$

that originates from the decomposition of each  $\mu_k$  and  $\Sigma_k$  into

$$\mu_k = \begin{pmatrix} \mu_k^{(s,a)} \\ \mu_k^Q \end{pmatrix} \quad \Sigma_k = \begin{pmatrix} \Sigma_k^{(s,a),(s,a)} & \Sigma_k^{(s,a),Q} \\ \Sigma_k^{Q,(s,a)} & \Sigma_k^{Q,Q} \end{pmatrix}, \quad (5)$$

where  $\mu_k(Q|s, a)$ ,  $\sigma_k^2(Q|s, a)$ ,  $y_g(s, a)$  and  $\beta_k(s, a)$  are respectively

$$\mu_k(Q|s, a) = \mu_k^Q + \Sigma_k^{Q,(s,a)} (\Sigma_k^{(s,a)})^{-1} ((s, a) - \mu_k^{(s,a)}) \quad (6)$$

$$\sigma_k^2(Q|s, a) = \Sigma_k^{Q,Q} - \Sigma_k^{Q,(s,a)} (\Sigma_k^{(s,a)})^{-1} \Sigma_k^{(s,a),Q} \quad (7)$$

$$y_g(s, a) = (\sigma_g^2(Q|s, a) + \mu_g^2(Q|s, a)) \quad (8)$$

$$\beta_k(s, a) = \frac{N(s, a; \mu_k^{(s,a)}, \Sigma_k^{(s,a),(s,a)})}{\sum_{k=1}^K N(s, a; \mu_k^{(s,a)}, \Sigma_k^{(s,a),(s,a)})}. \quad (9)$$

### B. Algorithm

$Q$ -CP is a planning algorithm that uses  $Q$ -values to effectively explore the state space  $S$  in a Monte-Carlo tree search, with the goal of generating action policies  $\pi_j$  for each agent  $j$  that operates in a team with a common goal. At each iteration  $i$ ,  $Q$ -CP(1) runs – for each agent – an UCT search, where admissible actions are selected through  $Q$ -value estimates, and (2) incrementally collects new samples that are used to improve  $Q$ -value estimates according to the procedure (as described in previous section).

More in detail (see Algorithm 1),  $Q$ -CP takes as input an initial policy  $\pi_j^0$  for each agent  $j$  and, at each iteration  $i = 1 \dots I$ , evolves as follows:

- 1) all the agents simultaneously follow their policy  $\pi_j^{i-1}$  for  $T$  timesteps, generating a set of  $T$  states  $\{s_t\}$ .
- 2) for each generated state  $s_t$ , each agent  $j$  sequentially runs a modified UCT search with depth  $H$ . Specifically, at each  $h = 1 \dots H$ , the search algorithm

- a) evaluates a subset of “admissible” actions  $\tilde{A}_j \subseteq A_j$  in  $s_{t+(h-1)}$ . Admissible actions are determined according to  $\hat{Q}(s_{t+(h-1)}, a_j)$  and  $\Sigma(Q|s_{t+(h-1)}, a_j)$ . In particular, differently from vanilla UCT, we only allow actions  $a_j$  such that
$$\hat{Q}(s_{t+(h-1)}, a_j) \geq \lambda \max_a \hat{Q}_j^{i-1}(s_{t+(h-1)}, a) - \delta$$

$$\delta \sim \sigma^2(Q|s_{t+(h-1)}, a_j), \quad (10)$$

where  $\lambda$  is typically initialized to 0.5 and increases with the number of iterations  $i$ . Through  $\delta$ , the prediction error for each action is captured, leading to a more directed exploration strategy. Additionally, an action can be randomly determined to be admissible with  $\varepsilon_{\tilde{A}}$  probability.

- b) selects and executes the best action  $a_j^{h*} \in \tilde{A}_j$  according to

$$\begin{aligned} e &= C \cdot \sqrt{\frac{\log(\sum_{a_j} \eta(s_{t+(h-1)}, a_j))}{\eta(s_{t+(h-1)}, a_j)}} \\ a_j^{h*} &= \arg \max_{a_j} \hat{Q}_j^{i-1}(s_{t+(h-1)}, a_j) + e, \end{aligned} \quad (11)$$

---

**Algorithm 1:  $Q$ -CP**

---

**Data:**  $n$  number of agents;  $I$  the number of iterations;  $\Delta$  initial state distribution;  $H$  UCT horizon;  $T$  policy execution timesteps,  $\lambda_0$  initial max.  $Q$  threshold multiplier for admissible actions,  $\varepsilon_{\bar{A}}$  probability for random admissible actions,  $\alpha$  learning rate,  $\gamma$  discount factor.

**Input:**  $\pi_j^0$  initial policy for each agent  $j$

**Output:**  $\pi_j^I$  policy learned after  $I$  iterations

**begin**

**for**  $i = 1$  **to**  $I$  **do**

$s_0 \leftarrow$  random state from  $\Delta$ .

**for**  $t = 1$  **to**  $T$  **do**

1)            Get state  $s_t$  by executing  $\pi_j^{i-1}(s_{t-1})$  for each agent  $j$ .

2)            **for**  $j = 1 \dots n$  **do**

3)                 $D_j^i \leftarrow \text{UCT}_{Q\text{-CP}}(s_t, \lambda_0, \varepsilon_{\bar{A}})$ .

4)                 $D_j^{0:i} \leftarrow D_j^i \cup D_j^{0:i-1}$ .

$\hat{Q}_j^i \leftarrow \text{function\_approximation}(D_j^{0:i}, \alpha, \gamma)$ .

$\pi_j^i(s) \leftarrow \arg \max_a \hat{Q}_j^i(s, a)$ .

**end**

**end**

**end**

**return**  $\pi_j^I$

**end**

---

where  $C$  is a constant that multiplies and controls the exploration term  $e$ , and  $\eta(s_{t+(h-1)}, a_j)$  is the number of occurrences of  $a_j$  in  $s_{t+(h-1)}$ . To use  $Q$ -CP on robotic applications, where the state space  $S$  is typically continuous, we define a comparison operator that introduces a discretization by informing the algorithm whether the difference of two states is smaller than a given threshold  $\xi$ .

- c) runs  $M$  simulations (or roll-outs), by simultaneously executing the  $\varepsilon$ -greedy policies of each agent  $j$  – based on  $\pi_j^{i-1}$  – until a terminal state is reached. It is crucially important to remark that during each roll-out, all the robots execute their policy. In this way,  $Q$ -CP evaluates the optimal action according teammates’ expected actions.

$Q$ -CP uses UCT as an expert and collects, at each iteration and for each agent, a dataset  $D_j^i$  of  $H$  samples  $x = (s, a_j, Q_j^i(s, a_j))$ .

- 3) after each iteration of UCT, each agent aggregates  $D_j^i$  into its own  $D_j^{0:i} = D_j^i \cup D_j^{0:i-1}$  [6], [28], and uses the complete dataset to learn the  $Q$ -values, as illustrated in previous section. It is worth remarking that selected actions influence the distribution of states and make the dataset non-i.i.d.
- 4) once  $\hat{Q}_j^i$  is learned, the policy is also updated as  $\pi_j^i(s) = \arg \max_a \hat{Q}_j^i(s, a)$ .

#### IV. EXPERIMENTAL EVALUATION

Generating competitive policies in a multi-robot environment is a challenging task, due to the exploration of large state spaces. Hence, in order to be practical in multiple

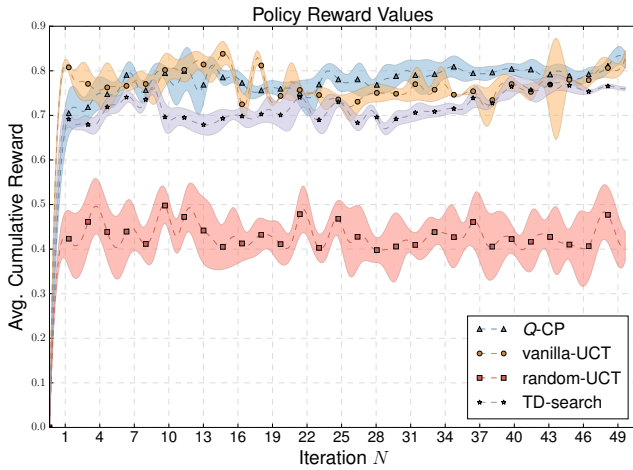
robotic applications, the goal of  $Q$ -CP is to enable a team of robots to plan an effective policy with a reduced number of explored states and algorithm iterations. Thus improving the simulation time and computational load. To this end, we compare our solution against TD-search [2] and both a *vanilla-UCT* and *random-UCT* implementations. We refer to *vanilla-UCT* as the standard UCT algorithm that, at each iteration, expands every possible action in  $A_j$ , for every agent  $j$ . The *random-UCT* instead, is a baseline algorithm, where at each step of the UCT search, it randomly selects only one action to expand. Specifically, we evaluate  $Q$ -CP in the three scenarios shown in Fig. 1: (1) a *cooperative navigation* application, where three robots have to coordinate in order to find the minimum path to their respective targets in a grid world; (2) a *hand-over task*, where two robots have to overcome structural constraints in order to complete their job; and (3) a *coordination task* among two robots passing through a door in opposite directions – this is a challenging setup, since timing and the selection of the correct action are critical.

#### A. Experimental Setup

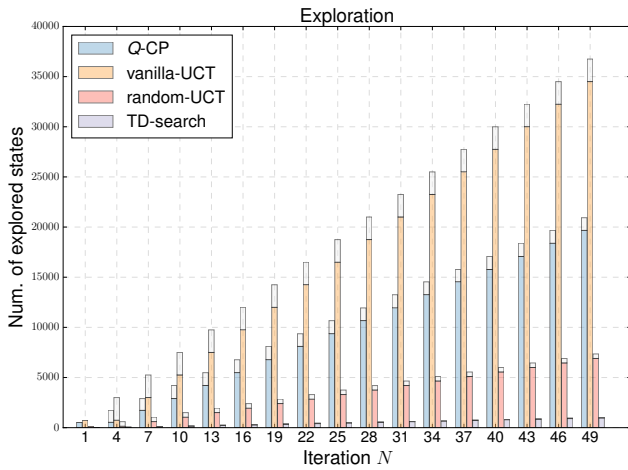
Experiments have been conducted using the V-REP simulator running on a single Intel Core i7-5700HQ core, with CPU@2.70GHz and 16GB of RAM. For all the scenarios, unless otherwise specified, the algorithm was configured with the same meta-parameters. The UCT horizon is set to  $H = 4$  to trade-off between usability and performance of the search algorithm; the number of roll-outs is set to be 3, while admissible actions are evaluated with an initial  $\lambda = 0.5$ , and  $\varepsilon_{\bar{A}} = 0.3$ , guaranteeing good amounts of exploration. The  $C$  constant in Eq. 11 is set to 0.7. The learning rate  $\alpha$  in Eq. 1 is set to 0.2, while the discount factor  $\gamma$  is equal to 0.8. The number of  $K$  components for the GMM approximator is 5 for the first two scenarios and 6 for *door passing*. The state space, the set of actions and the reward functions are finally implemented depending on the robots and applications. In each of the proposed applications, stochastic actions are obtained by randomizing the their outcomes with a 5% probability. We evaluate the reward obtained during different executions of  $Q$ -CP against the number of explored states and iteration of the algorithms. In particular, the learning system implements a *shaped* reward function, which computes a reward value for each of the visited states.

#### B. Cooperative Navigation

In this scenario, three robots have to perform a cooperative navigation task in a 4x4 grid world Fig. 1(a). Specifically, their task consists in reaching the square that matches their color (highlighted in the figure) by avoiding collisions with teammates. The state is a 12-features vector, where each agent is represented by 4 features encoding its current and target cells  $s = \langle r_x, r_y, t_x, t_y \rangle$ , while the set of actions, for each robot, is composed as a tuple of 5 elements  $A = \langle \text{noop}, \text{up}, \text{down}, \text{right}, \text{left} \rangle$ . The reward function is implemented to be inversely proportional to the sum



(a) Rewards



(b) States

Fig. 2. Normalized average reward (a) and number of explored states (b) obtained by  $Q$ -CP, TD-search,  $random$ -UCT and  $vanilla$ -UCT in 49 iterations in the cooperative navigation scenario. For each of them, the reward is averaged over 10 runs. The averaged reward is represented as a dotted line while, the line width represent its standard deviation. At each iteration, the number of explored states is represented as a bar, where the height indicates the total number of visited states, and the gray top highlights the amount of states added during the particular iteration  $i$ . The plot shows a comparison among  $Q$ -CP, TD-search,  $random$ -UCT and  $vanilla$ -UCT.

of the Manhattan distances among robots and their targets. Fig. 2 shows the comparison of  $Q$ -CP with TD-search,  $random$ -UCT,  $vanilla$ -UCT. Fig. 2(a) shows the normalized average reward (dotted line) and standard deviation (line width) of the considered algorithms obtained in 49 iterations and averaged over 10 runs for each of them. Fig. 2(b), shows the number of states explored per iteration. Each bar in the plot highlights the total number of states explored until the  $i$ -th iteration, and the gray top highlights the amount of states expanded during  $i$  with respect to the total number of states explored until  $i - 1$ . As expected,  $random$ -UCT explores a low number of states – since it expands only one random action at each UCT iteration – but it performs poorly. Conversely, if we consider  $Q$ -CP and  $vanilla$ -UCT, we

can notice similar reward values. Nonetheless, our solution generates a competitive policy with a remarkably reduced number of explored states. Although in this scenario TD-search explores less states than  $Q$ -CP, obtaining comparable rewards, the state space is relatively small (i.e. 4x4 grid world) and the problem can be solved by greedily following the estimated value function. In fact, once a first approximation is found at every state, following the max valued action matches the global optimum. As we will see in the following scenarios, however, TD-search reports a worse performance as the state dimensionality increases and, a greedy search over approximated value functions is no longer sufficient in both of them. In the case of the door passing experiment Fig. 1(b), in which immediately following the estimated value function – without allowing optimistic exploration – leads to deadlock in the robot policies. For example, if the goal is on the opposite side of the wall and the door is far from the robot, the agent needs to visit states that do not necessarily maximize a first approximation of the value function. Conversely, visiting these states (as  $Q$ -CP does by constructing upper confidence bounds of the value function), allows for a better approximation and thus for an improved policy.

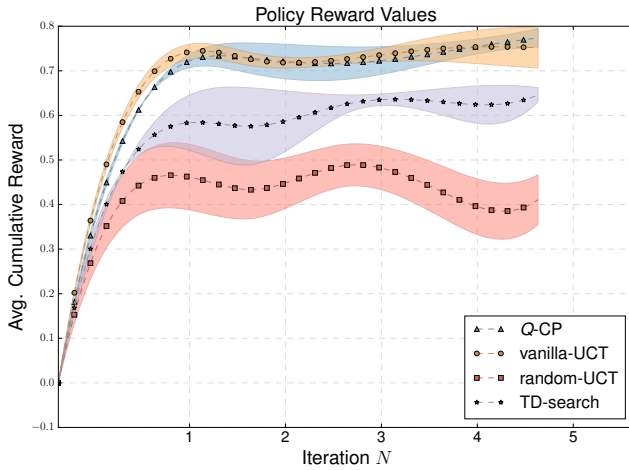
### C. Hand-over

In this scenario, two robots have to learn how to complete a hand-over task in order to move an object (Fig. 1(c)). Here, the state space is inherently continuous and is represented as a vector of 5 elements encoding the relative distance of both the robot bases (2D vector) and the end-effectors of their 5DOF arms (3D vector),  $s = \langle rel\_base_x, rel\_base_y, rel\_end\_effector_x, rel\_end\_effector_y, rel\_end\_effector_z \rangle$ . Note that, for running UCT, a discretization is performed as described in Section III-B. The set of 10 actions, for each robot, consists of  $A = \{ arm\text{-up}, arm\text{-down}, arm\text{-forward}, arm\text{-backward}, arm\text{-right}, arm\text{-left}, base\text{-forward}, base\text{-backward}, base\text{-left}, base\text{-right} \}$ . The reward function is implemented as a weighted sum of 2 components inversely proportional to the Euclidean distance between (1) the robot bases and a desired distance (1.0 m), and (2) the arm end-effectors. Hence, the reward promotes states where the two bases maintain a distance of one meter and their end-effectors are as close as possible to enable the hand-over of the object. As in the previous scenario, Fig. 3 shows a comparison between  $Q$ -CP, TD-search,  $random$ -UCT and  $vanilla$ -UCT. This scenario further highlights the ability of  $Q$ -CP in generating competitive policies with a reduced number of explored states. While TD-search and  $random$ -UCT report worse performance in terms of obtained reward,  $vanilla$ -UCT has more competitive values (Fig. 3(a)). However, the number of expanded nodes is significantly lower ( $< \sim 45\%$ ) since first iterations.

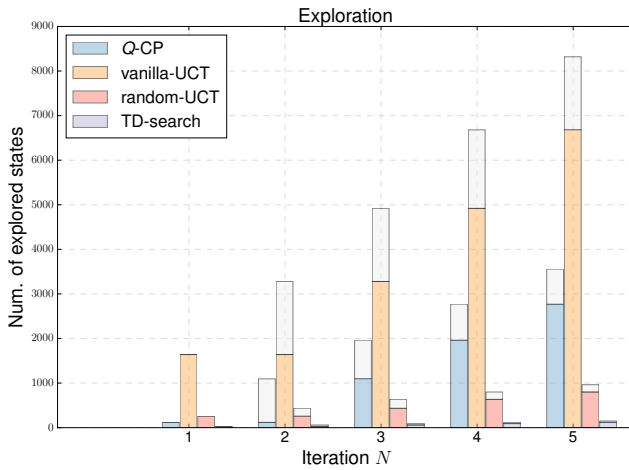
### D. Door Passing

The last scenario involves timing and precision in completing a coordination task, where two robots have to traverse a





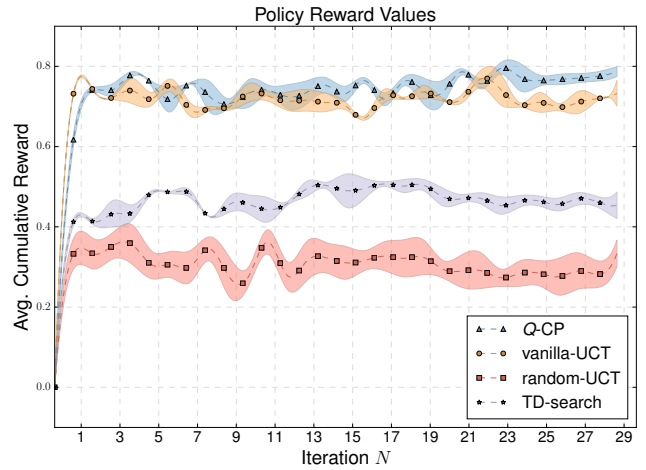
(a) Rewards



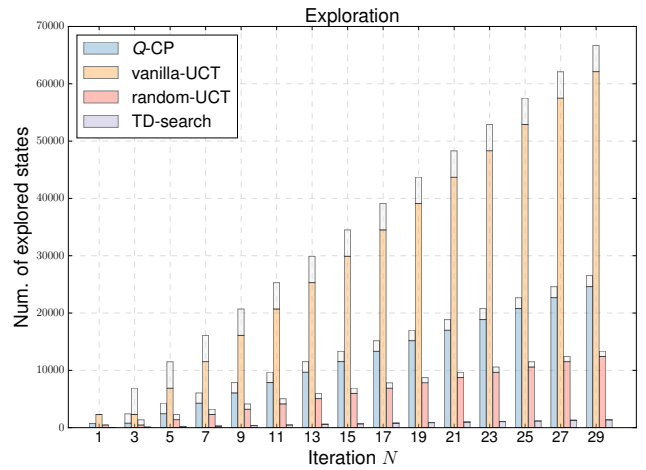
(b) States

Fig. 3. Normalized average reward and number of explored states obtained by *Q-CP*, *TD-search*, *random-UCT* and *vanilla-UCT* in 5 iterations in the hand-over scenario. For each of them, the reward is averaged over 10 runs.

narrow passage, e.g. a door Fig. 1(b). Also in this case, the problem is formalized by discretizing the continuous state space of the environment. Thus, the state is composed by an 8 features vector where each robot is represented by its position and the target one, as  $s = \langle r_x, r_y, t_x, t_y \rangle$ . Similarly to the first scenario, the robots share a set of 5 actions being  $A = \langle \text{noop}, \text{up}, \text{down}, \text{right}, \text{left} \rangle$ . The reward function is a weighted sum of 3 components: (1) the first component is inversely proportional to the distance among robots and their targets, (2) the second component is proportional to the distance to obstacles and induces the generated policy to prefer paths away from them, while (3) the third component is proportional to the distance between the two robots penalizing states where they are too close. Fig. 4 evaluates *Q-CP*, *TD-search*, *random-UCT* and *vanilla-UCT* over 29 iterations. The analysis of the plots is similar to the previous scenarios and highlights that *Q-CP* presents the most suitable trade-off between performance and computational load, and thus results to be more practical



(a) Rewards



(b) States

Fig. 4. Normalized average reward and number of explored states obtained by *Q-CP*, *TD-search*, *random-UCT* and *vanilla-UCT* in 29 iterations in the cooperative navigation scenario. For each of them, the reward is averaged over 10 runs.

in robotics applications. In fact, with respect to *vanilla-UCT*, the number of explored states is still reduced by the 60% during last iterations. *TD-search*, conversely, is not able to generate a competitive policy by greedily following the approximated value function (as explained in Sec IV-B). This gives a substantial advantage to *Q-CP* in planning complex tasks as the one in this scenario. It is worth remarking that these plots demonstrate the ability of *Q-CP* to guide the learning routine towards promising areas of the state search since the first iteration. In fact, *Q-CP* is able to (1) crop out states that do not match the objective of the task and to (2) generate performing policies even with a reduced number of samples and iterations.

## V. CONCLUSION

In this paper we introduced *Q-CP*, an iterative policy generation algorithm that uses action values to guide and reduce the exploration of the state space in different multi-

robot scenarios. Our key contribution is the combination of a Monte-Carlo tree search algorithm with action values, which demonstrates a remarkable improvement in the computational load of the algorithm without loss in the performance. Such an improvement makes our approach more practical and suitable in difficult robotics applications, where the lack of training examples is often a limiting condition.

$Q$ -CP presents several limitations that give the opportunity to further improve our solution. The major problem is the massive use of simulation calls, that make the algorithm less appealing. Hence, we are working on an online version by relying on online and offline knowledge as in [29]. Moreover, as the task complexity increases, the non-linearity of the value function is a structural problem that affects performance. Thus, we aim at exploiting different function approximations in order to improve  $Q$ -CP and apply it with larger state spaces. For example, based on several recent advances in policy learning [30], we are evaluating the use of neural networks or deep networks to further improve the performance of our approach by defining more informed state representations and function approximators.

#### REFERENCES

- [1] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of monte carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, no. 1, pp. 1–43, 2012.
- [2] D. Silver, R. S. Sutton, and M. Müller, "Temporal-difference search in computer go," *Machine Learning*, vol. 87, no. 2, pp. 183–219, 2012.
- [3] L. Panait and S. Luke, "Cooperative multi-agent learning: The state of the art," *Autonomous Agents and Multi-Agent Systems*, vol. 11, no. 3, pp. 387–434, Nov. 2005. [Online]. Available: <http://dx.doi.org/10.1007/s10458-005-2631-2>
- [4] L. Kocsis and C. Szepesvári, "Bandit based monte-carlo planning," *Machine Learning: ECML 2006*, pp. 282–293, 2006.
- [5] A. Agostini and E. Celaya, "Reinforcement learning with a gaussian mixture model," in *The 2010 International Joint Conference on Neural Networks*, July 2010, pp. 1–8.
- [6] S. Ross, G. J. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *International Conference on Artificial Intelligence and Statistics*, 2011, pp. 627–635.
- [7] M. H. Winands, Y. Björnsson, and J.-T. Saito, "Monte-carlo tree search solver," in *Proceedings of the 6th International Conference on Computers and Games*, ser. CG '08. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 25–36. [Online]. Available: [http://dx.doi.org/10.1007/978-3-540-87608-3\\_3](http://dx.doi.org/10.1007/978-3-540-87608-3_3)
- [8] M. Bowling, "Convergence problems of general-sum multiagent reinforcement learning," in *In Proceedings of the Seventeenth International Conference on Machine Learning*. Morgan Kaufmann, 2000, pp. 89–94.
- [9] C. Claus and C. Boutilier, "The dynamics of reinforcement learning in cooperative multiagent systems," in *Proceedings of the Fifteenth National/Tenth Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence*, ser. AAAI '98/IAAI '98. Menlo Park, CA, USA: American Association for Artificial Intelligence, 1998, pp. 746–752. [Online]. Available: <http://dl.acm.org/citation.cfm?id=295240.295800>
- [10] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *International Journal of Robotics Research*, July 2013.
- [11] S. Nikolaidis, R. Ramakrishnan, K. Gu, and J. Shah, "Efficient model learning from joint-action demonstrations for human-robot collaborative tasks," in *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction*, ser. HRI '15. ACM, 2015, pp. 189–196.
- [12] B. Kim and J. Pineau, "Socially adaptive path planning in human environments using inverse reinforcement learning," *International Journal of Social Robotics*, vol. 8, no. 1, pp. 51–66, 2016.
- [13] R. S. Sutton, D. Precup, and S. Singh, "Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning," *Artificial Intelligence*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [14] G. Konidaris, "Constructing abstraction hierarchies using a skill-symbol loop," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, ser. IJCAI'16. AAAI Press, 2016, pp. 1648–1654. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3060832.3060851>
- [15] T. G. Dietterich, "Hierarchical reinforcement learning with the maxq value function decomposition," *Journal of Artificial Intelligence Research (JAIR)*, vol. 13, pp. 227–303, 2000.
- [16] R. Makar, S. Mahadevan, and M. Ghavamzadeh, "Hierarchical multi-agent reinforcement learning," in *Proceedings of the fifth international conference on Autonomous agents*. ACM, 2001, pp. 246–253.
- [17] F. Abtahi and M. R. Meybodi, "Solving multi-agent markov decision processes using learning automata," in *2008 6th International Symposium on Intelligent Systems and Informatics*, Sept 2008, pp. 1–6.
- [18] D. Hadfield-Menell, A. D. Dragan, P. Abbeel, and S. J. Russell, "Cooperative inverse reinforcement learning," *CoRR*, vol. abs/1606.03137, 2016. [Online]. Available: <http://arxiv.org/abs/1606.03137>
- [19] D. S. Bernstein, C. Amato, E. A. Hansen, and S. Zilberstein, "Policy iteration for decentralized control of markov decision processes," *CoRR*, vol. abs/1401.3460, 2014. [Online]. Available: <http://arxiv.org/abs/1401.3460>
- [20] S. Proper and P. Tadepalli, "Solving multiagent assignment markov decision processes," in *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems - Volume 1*, ser. AAMAS '09. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2009, pp. 681–688.
- [21] Y. Engel, S. Mannor, and R. Meir, "Reinforcement learning with gaussian processes," in *Proceedings of the 22Nd International Conference on Machine Learning*, ser. ICML '05. New York, NY, USA: ACM, 2005, pp. 201–208. [Online]. Available: <http://doi.acm.org/10.1145/1102351.1102377>
- [22] G. Chowdhary, M. Liu, R. Grande, T. Walsh, J. How, and L. Carin, "Off-policy reinforcement learning with gaussian processes," *IEEE/CAA Journal of Automatica Sinica*, vol. 1, no. 3, pp. 227–238, 2014.
- [23] S. Gelly and D. Silver, "Monte-carlo tree search and rapid action value estimation in computer go," *Artif. Intell.*, vol. 175, no. 11, pp. 1856–1875, July 2011. [Online]. Available: <http://dx.doi.org/10.1016/j.artint.2011.03.007>
- [24] S. James, G. Konidaris, and B. Rosman, "An analysis of monte carlo tree search," in *AAAI*, 2017, pp. 3576–3582.
- [25] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977.
- [26] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability - Vol. 1*, L. M. Le Cam and J. Neyman, Eds. University of California Press, Berkeley, CA, USA, 1967, pp. 281–297.
- [27] G. Schwarz *et al.*, "Estimating the dimension of a model," *The annals of statistics*, vol. 6, no. 2, pp. 461–464, 1978.
- [28] S. Ross and J. A. Bagnell, "Reinforcement and imitation learning via interactive no-regret learning," *arXiv preprint arXiv:1406.5979*, 2014.
- [29] S. Gelly and D. Silver, "Combining online and offline knowledge in uct," in *Proceedings of the 24th international conference on Machine Learning*. ACM, 2007, pp. 273–280.
- [30] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 02 2015. [Online]. Available: <http://dx.doi.org/10.1038/nature14236>