

Simultaneous Discovery of Multiple Alternative Optimal Policies by Reinforcement Learning

Petar Kormushev and Darwin G. Caldwell

Department of Advanced Robotics

Istituto Italiano di Tecnologia

Via Morego 30, 16163 Genova, Italy

Email: {petar.kormushev, darwin.caldwell}@iit.it

Abstract—Conventional reinforcement learning algorithms for direct policy search are limited to finding only a single optimal policy. This is caused by their local-search nature, which allows them to converge only to a single local optimum in policy space, and makes them heavily dependent on the policy initialization.

In this paper, we propose a novel reinforcement learning algorithm for direct policy search, which is capable of simultaneously finding multiple alternative optimal policies. The algorithm is based on particle filtering and performs global search in policy space, therefore eliminating the dependency on the policy initialization, and having the ability to find the globally optimal policy. We validate the approach on one- and two-dimensional problems with multiple optima, and compare its performance to a global random sampling method, and a state-of-the-art Expectation-Maximization based reinforcement learning algorithm.

Keywords—reinforcement learning; particle filters; global search

I. INTRODUCTION

Reinforcement learning (RL) is a machine learning approach, in which the goal is to find a policy π that maximizes the expected future return, calculated based on a scalar reward function $R(\cdot) \in \mathbb{R}$. The argument of $R(\cdot)$ can be defined in different ways, e.g. it could be a state s , or a state transition, or a state-action pair, or a whole trial as in the case of episodic RL, etc. The policy π determines what actions will be performed by the RL agent, and is usually state dependent. [1]

Originally, the RL problem was formulated in terms of a Markov Decision Process (MDP) or Partially Observable MDP (POMDP). In this formulation, the policy π is viewed as a direct mapping function ($\pi : s \mapsto a$) from state $s \in S$ to action $a \in A$.

Alternatively, instead of trying to learn the explicit mapping from states to actions, it is possible to perform *direct policy search*, as shown in [2]. In this case, the policy π is considered to depend on some parameters $\theta \in \mathbb{R}^N$, and is written as a parameterized function $\pi(\theta)$. The episodic reward function becomes $R(\tau(\pi(\theta)))$, where τ is a trial performed by following the policy. The reward can be abbreviated as $R(\tau(\theta))$ or even as $R(\theta)$, which reflects the idea that the behaviour of the RL agent can be influenced by only changing the values of the policy parameters θ . Therefore, the outcome of the behaviour, which is represented by the reward $R(\theta)$, can be optimized by only optimizing the values θ . This way, the RL problem is

transformed into a black-box optimization problem with cost function $R(\theta)$, as shown in [3] under the name *parameter-based exploration*.

However, it is infeasible to use conventional numeric optimization methods to maximize $R(\theta)$ if we want to apply RL to real-world problems, because the cost function is usually expensive to evaluate. For example, each cost function evaluation requires conducting at least one trial which could involve costly real-world experiments or computationally expensive simulations. Therefore, alternative optimization methods are desired, which are tuned to the specific need of RL to reduce as much as possible the number of reward function evaluations (trials).

This paper aims to bring new ideas into the domain of RL, borrowed from statistics. We propose a novel direct policy search RL algorithm which provides an alternative solution to the RL problem, and new possibilities for RL algorithms in general. Before we can introduce our proposed approach, we need to place it in the proper context, which is done in the following Sections II and III.

II. STATE-OF-THE-ART RL ALGORITHMS FOR DIRECT POLICY SEARCH

This section contains a non-exhaustive list of direct policy search RL approaches. We focus on policy search methods that attempt to minimize the number of trials. Examples for such RL methods include:

- Policy Gradient based RL - in which the RL algorithm is trying to estimate the gradient of the policy with respect to the policy parameters, and to perform gradient descent in policy space. The Episodic Natural Actor-Critic (eNAC), in [4], and Episodic REINFORCE in [5] are two of the well-established approaches of this type.
- Expectation-Maximization based RL - in which the EM algorithm is used to derive an update rule for the policy parameters at each step, trying to maximize the lower bound on the expected return of the policy. A state-of-the-art RL algorithm of this type is PoWER (Policy learning by Weighting Exploration with the Returns), in [6], as well as its generalization MCEM, in [7].
- Path Integral based RL - in which the learning of the policy parameters is based on the framework of stochastic optimal control with path integrals. A state-of-the-art RL

algorithm of this type is PI² (Policy Improvement with Path Integrals), in [8].

- Regression based RL - in which regression is used to calculate updates to the RL policy parameters using the rewards as weights. One of the approaches that are used extensively is LWPR (Locally Weighted Projection Regression), in [9].
- Model-based policy search RL - in which a model of the transition dynamics is learned and used for long-term planning. Policy gradients are computed analytically for policy improvement using approximate inference. A state-of-the-art RL algorithm of this type is PILCO (Probabilistic Inference for Learning COntrol), in [10].

A major problem with all these existing approaches is that they only perform *local search*, therefore they do not guarantee convergence to the global optimum. Instead, all of them tend to converge to some local sub-optimal solution.

Another major problem is that the result from each method is largely influenced by the initial value of θ . The reason for this is that all of these methods have the inherent notion of ‘current policy’. They work by changing this current policy iteratively with a small amount every time using an update rule in the following generic form:

$$\theta_{n+1} = \theta_n + \Delta\theta, \quad (1)$$

where θ_n is the current policy and θ_{n+1} is the new policy. The term $\Delta\theta$ is calculated in different ways, either by some form of gradient estimation (as in gradient-based RL algorithms and PILCO), or using other calculation method (e.g. EM based, path integral based or regression-based). Regardless of the exact way, these approaches are, by definition, local search methods, and they can only guarantee convergence to a local optima.

In contrast, what we propose in this paper is a *global search* method for direct policy search reinforcement learning which is guaranteed not to get stuck at local optima. The principle for realizing this comes from another research domain: *particle filtering*.

III. PARTICLE FILTERS

Particle filters, also known as Sequential Monte Carlo methods [11], [12], originally come from statistics and are similar to importance sampling methods. Particle filters are able to approximate any probability density function, and can be viewed as a ‘sequential analogue’ of Markov chain Monte Carlo (MCMC) batch methods.

Although particle filters are mainly used in statistics, there are a few other research areas in which particle filters have found application. For example, in the domain of probabilistic robotics [13], particle filters are extensively and successfully used, e.g. for performing Monte Carlo localization of mobile robots with respect to a global map of the terrain [14], [15], and also for Simultaneous Localization and Mapping (SLAM) tasks [16], [17], [18].

The potential of applying particle filters in the RL domain appears to be largely unexplored so far. To the best of our

knowledge, there are only two partial attempts to apply particle filters in RL in the existing published work, done by Notsu et al and Samejima et al, respectively, as follows.

In [19], they studied traditional RL with discrete state and action spaces. They used the Actor-Critic method for performing value iteration, with state value updates using the TD-error. In this conventional framework, they proposed a particle filter for segmentation of the action space. Their goal was to do domain reduction, i.e. to minimize the number of discrete actions available at each state, by dividing the action space in segments. Then, in [20], they extended the same idea to traditional continuous RL and used Q-learning with function approximation. They used a similar particle filter approach for segmenting the continuous state and action spaces into discrete sets by particles, and applied it to inverted pendulum tasks.

In [21], they studied neurophysiology and created a reinforcement learning model of an animal behaviour. Their goal was to predict the behaviour of a monkey during an experimental task. They used traditional Q-learning, and built a Bayesian network representation of the Q-learning agent. In this framework, particle filtering was used to estimate action probability in order to predict the animal behaviour.

In both of these existing approaches, particle filters were used in a limited way, as a technique to solve some partial problem within a traditional RL framework.

In this paper, we propose a rather different approach. First, we propose a completely new view of the link between particle filters and RL. Then, we propose an entirely novel RL algorithm for direct global policy search, based on particle filters as the core for the RL algorithm itself. In our framework, the search is performed in the policy space defined by the selected policy parameterization, and the process is viewed as a black-box optimization. The particle filter itself is the core of the proposed RL algorithm, and is responsible for guiding the exploration and exploitation, by creating particles, each of which represents a whole policy. Details of the proposed novel view and algorithm follow.

IV. NOVEL VIEW OF RL AND ITS LINK TO PARTICLE FILTERS

The key to linking particle filters and RL is to make the following observation. The landscape, defined by the reward function $R(\theta) \in \mathbb{R}$ over the whole continuous domain of the parameter space $\theta \in \Theta$, can be viewed as defining an improper probability density function (IPDF)¹. This is possible even if the reward function $R(\theta)$ has negative values in its range, because we can simply add a constant positive number $L = |\inf_{\theta \in \Theta} R(\theta)|$ to it, and obtain a new reward function $R'(\theta)$ which is non-negative and has exactly the same set of optimizers $\theta^* \in \Theta$ as the original one. Therefore, optimizing $R'(\theta)$ will also optimize $R(\theta)$.

Once we make the assumption that $R(\theta)$ is just an IPDF, then the RL problem can be reformulated from a new point

¹IPDF is similar to PDF except that the integral of it does not have to be equal to one.

of view. Each trial $\tau(\pi(\theta))$ can be viewed as an independent sample from this unknown IPDF. The RL algorithm can be viewed as a method for choosing a finite number of sampling points for which to obtain the value of the IPDF. Finally, the RL problem can be viewed as the problem of finding the mode (or all modes, in the multi-modal case) of the unknown IPDF, given only a finite number of sampling points with their corresponding values of the IPDF, obtained by the RL algorithm.

This view of RL immediately opens the path for applying particle filters, because they are a method for approximate estimation of an unknown PDF based on a finite number of samples. To complete the link between RL and particle filters, the only thing left to state is that it is trivial to convert an IPDF into a PDF simply by normalizing it (i.e. dividing by the integral of it).

V. RL BASED ON PARTICLE FILTERS

Using the reformulation of RL from the previous section, here we propose a novel RL algorithm based on Particle Filters (RLPF). The main idea of RLPF is to use particle filtering as a method for choosing the sampling points, i.e. for calculating a parameter vector θ for each trial.

We define a *policy particle* p_i to be the tuple $p_i = \langle \theta_i, \tau_i, R_i, w_i \rangle$, where the particle p_i represents the outcome of a single trial τ_i performed by executing an RL policy $\pi(\theta_i)$, where θ_i is a vector of policy parameter values modulating the behaviour of the RL policy π . The policy particle also stores the value of the reward function evaluated for this trial $R_i = R(\tau_i(\pi(\theta_i)))$. The variable τ_i contains task-specific information recorded during the trial depending on the nature of the task. The information in τ_i is used by the reward function to perform its evaluation. The variable w_i is the importance weight of this policy particle, and the way of its calculation is explained below.

Following closely the ideas from particle filters, we make the assumption that the set of particles $\{p_i\}$ is an approximate implicit representation of the underlying unknown IPDF defined by $R(\theta)$. Therefore, in order to select a new particle which obeys the real IPDF distribution, what we can do is to sample from the approximate distribution while correcting for the discrepancy between it and the real distribution. The mechanism for this correction is provided by the importance weights $\{w_i\}$.

Firstly, each policy particle p_i is assigned a scalar importance weight w_i derived from its corresponding reward R_i using a transformation function g , such that: $w_i \propto g(R_i)$. In the simplest case, $g(\cdot)$ could be the identity, but in the general case, it could be an arbitrary non-negative function. We apply the function g in such a way, that the importance weights are normalized, in the sense that: $\forall w_i \quad 0 < w_i < 1$, and also: $\sum w_i = 1$.

Secondly, we construct an auxiliary function $h(u) = \int_{-\infty}^u w_u du$, which in our discrete case takes the form $h(k) = \sum_{j=1}^k w_j$. This function can be thought of as the (approximate) cumulative density function (CDF) of the unknown PDF.

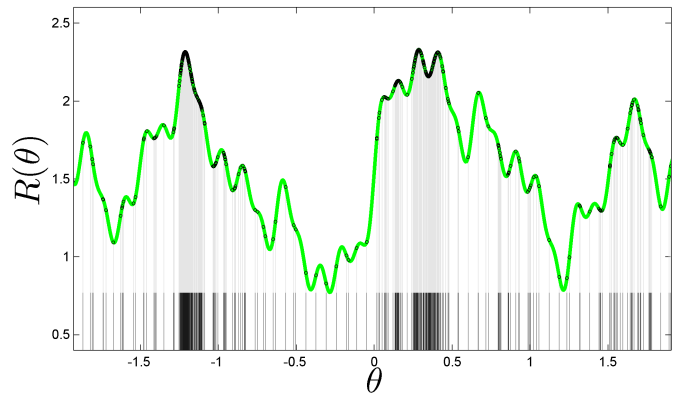


Fig. 1. An illustration of a typical run of RLPF on a 1-dimensional problem. The generated policy particles by RLPF are shown with vertical grey stripes. The corresponding reward values are shown with black circles on top of the reward function line, shown in green. The following synthetic reward function was used, because it has many local optima: $R(\theta) = 1.55 + 0.3 \sin(2.7\theta) + 0.4 \sin(4.5\theta) + 0.2 \sin(8.7\theta) + 0.14 \sin(20\theta) + 0.08 \sin(30\theta) + 0.08 \sin(50\theta)$.

Indeed, due to the way we create the importance weights, it follows directly that $\int_{-\infty}^{+\infty} w_u du = 1$, and thus $h(u)$ is a proper CDF. This is important because, given that $w_i > 0$, it guarantees that $h(u)$ is strictly monotonically increasing and therefore the inverse function h^{-1} exists.

Thirdly, we introduce a random variable z which is uniformly distributed in the interval $(0, 1)$. Now, it can be shown that the random variable y defined as $y = h^{-1}(z)$ is distributed (approximately) according to the desired unknown PDF, see e.g. [22].

At this point, there are two variants of particle filters, which accordingly result in two variants for RLPF, and they are:

- Sequential Importance Resampling (SIR) - this is the original particle filtering algorithm [23]. It approximates the filtering distribution by a weighted set of particles, very similar to what we described so far. However, unlike RLPF, where our goal is to find the mode (or all modes) of the unknown IPDF, the goal of SIR is to find the expectation of a function by approximating it as a weighted average.
- Sequential Importance Sampling (SIS) - This variant is very similar to SIR, except that it does not use the resampling stage.

In SIR, resampling is used in order to avoid the problem of degeneracy of the algorithm, that is, avoiding the situation that all but one of the importance weights are close to zero. The performance of the algorithm can be also affected by proper choice of resampling method.

However, the goal of RLPF is not to approximate the expectation of a function, but rather, to find the mode (or modes) of the unknown function $R(\theta)$. Thus, we describe the proposed RLPF algorithm as a special case of SIS, while keeping in mind that it can easily be extended into SIR by adding a resampling step. The pseudo-code for RLPF is given in Algorithm 1.

VI. ANALYSIS OF RLPF

RLPF inherits many advantages from particle filters, among which:

- It is very simple to implement and can be realized easily in embedded systems for online learning;
- It is not computationally expensive and has very low memory footprint;
- It can use adaptive computation depending on the available resources (both time- and CPU-wise) by changing the value of the σ parameter;
- It can concentrate the effort of the RL exploration on the most important parts of the policy space, by using function $g(R)$ which increases the relative difference between the rewards, e.g. the function $g(R) = (1 + R)^2$;
- It can exhibit adaptive convergence rate depending on the requirements for precision and time, by changing the initial noise level ϵ_0 and the decay factor λ .

Algorithm 1 Reinforcement Learning based on Particle Filters (RLPF)

Input: parameterized policy π , policy parameter space Θ , reward function $R(\theta)$ where $\theta \in \Theta$, reward transformation function g , total number of trials N , initialization number of trials $L < N$, initial noise ϵ_0 , noise decay factor λ , maximum number of particles σ .

```

Let  $S = \{ \}$  {A set of policy particles}
for  $l = 1$  to  $L$  do
  Draw  $\theta_l \sim U(\Theta)$  {Sample  $L$  initial particles}
  Perform trial  $\tau_l(\pi(\theta_l))$ 
  Create new policy particle  $p_l = \langle \theta_l, \tau_l, R_l, w_l \rangle$ 
   $S = S \cup \{p_l\}$ 
end for
for  $n = L + 1$  to  $N$  do
  Let  $h(0) = 0$ 
  for  $i = 1$  to  $|S|$  do
     $w_i = \frac{g(R_i)}{\sum_{j=1}^{|S|} g(R_j)}$  {Calc. importance weights}
     $h(i) = h(i-1) + w_i$  {Calc. aux. function}
  end for
  Draw  $z \sim U(0, 1)$ 
  Let  $y = h^{-1}(z)$ 
  Let  $k = \lceil y \rceil$  { $\lceil \cdot \rceil$  is the ceiling function}
  Select policy particle  $p_k = \langle \theta_k, \tau_k, R_k, w_k \rangle$ 
  Let  $\epsilon_n = \epsilon_0 \lambda^{(n-L-1)}$  {noise with exp. decay  $\lambda$ }
  Let  $\theta_n = \theta_k + \epsilon_n$ 
  Perform trial  $\tau_n(\pi(\theta_n))$ 
  Create new policy particle  $p_n = \langle \theta_n, \tau_n, R_n, w_n \rangle$ 
   $S = S \cup \{p_n\}$ 
  if  $|S| > \sigma$  then
    {Remove policy particle with smallest reward}
     $j = \operatorname{argmin}_{p_j \in S} g(R_j)$ 
     $S = S \setminus \{p_j\}$ 
  end if
end for

```

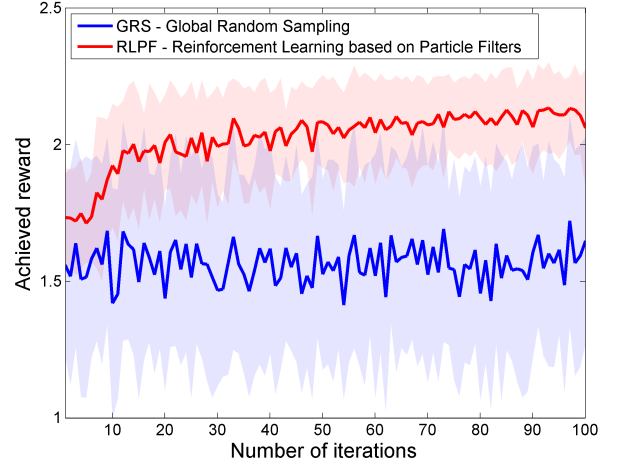


Fig. 2. A comparison of the convergence of two global policy search RL algorithms: Global Random Sampling (GRS) policy search RL algorithm, and RLPF. The results are averaged over 50 runs of each algorithm. Every run has 100 trials. RLPF easily outperforms GRS both in terms of achieved reward and low level of variance.

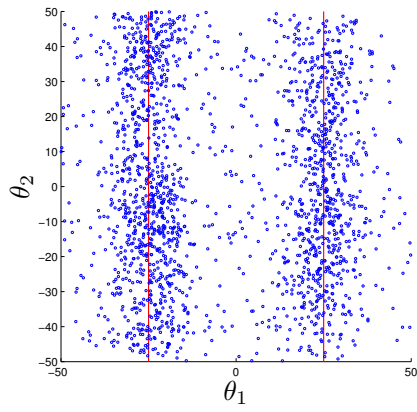
To be fair, we should also mention some disadvantages of RLPF. For example, as a member of global search methods, it generally requires more trials in order to converge, because the scope of the search is the largest possible - the whole policy space. However, using appropriate reward transformation functions, it is possible to bias the exploration towards the highest-reward particles, sacrificing thoroughness for convergence speed.

VII. EXPERIMENTAL EVALUATION OF RLPF

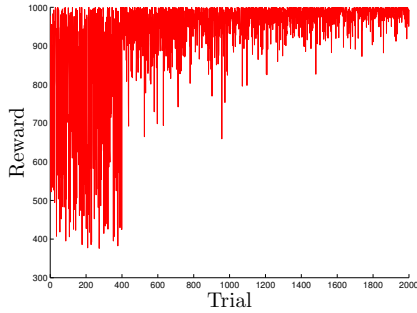
First, we evaluate RLPF standalone on a one-dimensional problem, because it is the easiest to visualize and analyze. Figure 1 shows an example run of RLPF on a class of synthetic 1D reward functions with many local optima. It is clearly visible that the generated policy particles by RLPF tend to cluster around the highest ‘peaks’ in the reward function.

Second, we compare the performance of RLPF with other global policy search RL methods. It is difficult to select baseline against which to compare RLPF to, because there are no any truly global search policy-based RL algorithms. It would not be fair to compare a local search RL, such as policy gradient based RL, to RLPF, because the local search methods will easily get stuck at the local optima. So, instead, for a baseline we use a stochastic global policy search RL algorithm, which is based on Global Random Sampling (GRS) in policy space. The comparison with RLPF is shown in Figure 2, averaged over many runs.

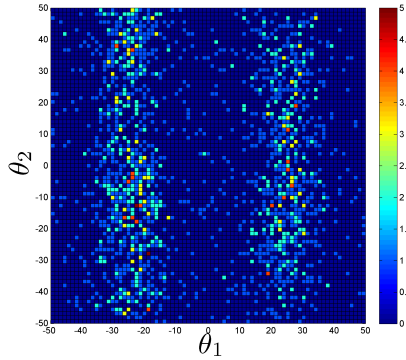
Third, we evaluate RLPF on a two-dimensional problem. Figure 3 shows an example run of RLPF on a 2D reward function, which has multiple alternative optima all having the same value. RLPF is able to efficiently explore the whole policy space, by covering it uniformly with policy particles during the initialization phase. Then, for the rest of the trials,



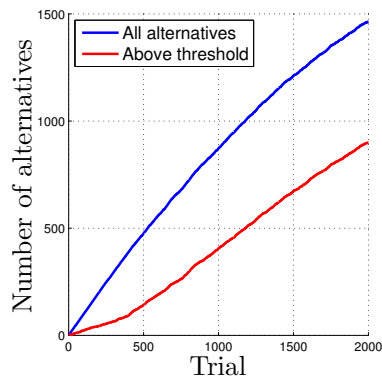
(a) Explored policies by RLPF in policy space



(b) Reward achieved at each trial by RLPF

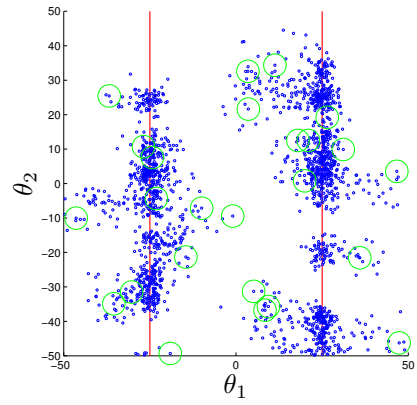


(c) 2D histogram of alternatives discovered by RLPF

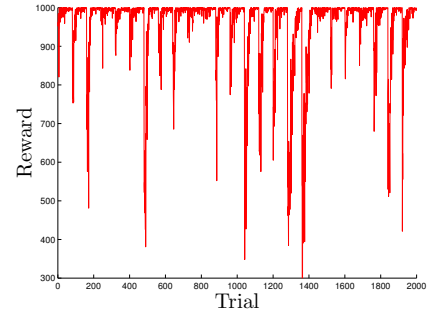


(d) Number of alternatives found vs. trials

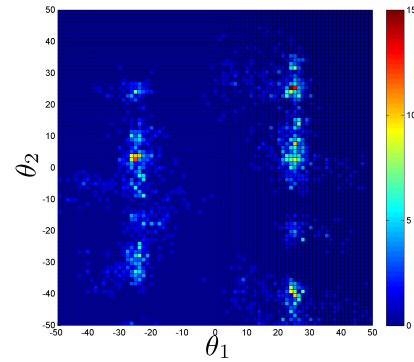
Fig. 3. Evaluation of the proposed RLPF algorithm on a 2D problem with the following reward function: $R(\theta) = 1000 - \min(|\theta_1 + 25|, |\theta_1 - 25|)^2$. The optimal policies lie on two vertical lines ($\theta_1 = \pm 25$), and are illustrated with red lines. The total number of trials is 2000 (the first 400 of them are random), and RLPF manages to find almost 1500 alternatives in total, with around 900 of them above the threshold (95% of the maximum reward).



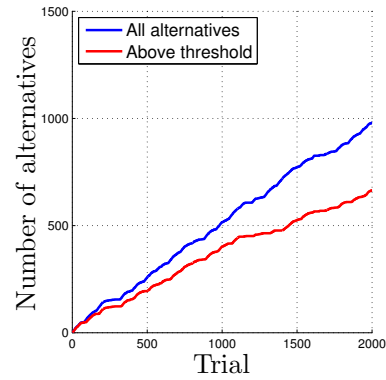
(a) Explored policies by PoWER in policy space



(b) Reward achieved at each trial by PoWER



(c) 2D histogram of alternatives discovered by PoWER



(d) Number of alternatives found vs. trials

Fig. 4. Comparative evaluation of PoWER algorithm on the same 2D problem as in Fig. 3. PoWER is run 25 times in batch mode for 80 trials for a total of 2000 trials (same as RLPF). Each run is started with a random initial policy. Due to the less efficient policy space exploration compared to RLPF, PoWER manages to find only less than 1000 alternatives in total, with around 700 of them above the threshold.

RLPF focuses the exploration on the policy space regions with highest rewards.

Fourth, we compare RLPF to a state-of-the-art policy-search RL algorithm. We chose the PoWER algorithm, because of its fast convergence, and small number of open parameters that need to be tuned. Since PoWER is a local-search RL method, in order to be fair in the comparison, we run it multiple times starting from a random initial policy every time. Figure 4 shows an example run of PoWER on the same 2D reward function as in the previous experiment with RLPF. Comparing the two figures shows the significant advantage of RLPF over PoWER both in terms of exploration and speed of discovery of alternative optimal policies.

VIII. CONCLUSION

This paper introduced ideas from particle filtering and importance sampling into the domain of reinforcement learning. We revealed a link between particle filters and RL, and reformulated the RL problem from this perspective. We proposed a new RL algorithm which is based on particle filters at its core. Due to the ability of particle filters to perform global search, the resulting RL algorithm is also capable of direct global search in policy space, which is significant improvement over traditional local search based policy RL. Since this work opens up a novel research direction in RL, there are many ways in which it can be extended in the future.

ACKNOWLEDGMENT

This work is partially supported by the European project PANDORA under contract FP7-ICT-288273.

REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Reinforcement learning: an introduction*, ser. Adaptive computation and machine learning. Cambridge, MA, USA: MIT Press, 1998.
- [2] M. Rosenstein and A. Barto, "Robot weightlifting by direct policy search," in *International Joint Conference on Artificial Intelligence*, vol. 17, no. 1. Citeseer, 2001, pp. 839–846.
- [3] T. Rückstieß, F. Sehnke, T. Schaul, D. Wierstra, Y. Sun, and J. Schmidhuber, "Exploring parameter space in reinforcement learning," *Paladyn. Journal of Behavioral Robotics*, vol. 1, no. 1, pp. 14–24, 2010.
- [4] J. Peters and S. Schaal, "Natural actor-critic," *Neurocomput.*, vol. 71, no. 7-9, pp. 1180–1190, 2008.
- [5] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Mach. Learn.*, vol. 8, no. 3-4, pp. 229–256, 1992.
- [6] J. Kober and J. Peters, "Learning motor primitives for robotics," in *Proc. IEEE Intl Conf. on Robotics and Automation (ICRA)*, May 2009, pp. 2112–2118.
- [7] N. Vlassis, M. Toussaint, G. Kontes, and S. Piperidis, "Learning model-free robot control by a Monte Carlo EM algorithm," *Autonomous Robots*, vol. 27, no. 2, pp. 123–130, 2009.
- [8] E. Theodorou, J. Buchli, and S. Schaal, "A Generalized Path Integral Control Approach to Reinforcement Learning," *The Journal of Machine Learning Research*, vol. 11, pp. 3137–3181, December 2010.
- [9] S. Vijayakumar and S. Schaal, "Locally weighted projection regression: An O(n) algorithm for incremental real time learning in high dimensional spaces," in *Proc. Intl Conf. on Machine Learning (ICML)*, Haifa, Israel, 2000, pp. 288–293.
- [10] 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 Machine Learning*, L. Getoor and T. Scheffer, Eds., Bellevue, WA, USA, June 2011.
- [11] A. Doucet, S. Godsill, and C. Andrieu, "On sequential monte carlo sampling methods for bayesian filtering," *Statistics and computing*, vol. 10, no. 3, pp. 197–208, 2000.
- [12] A. Doucet, N. De Freitas, and N. Gordon, *Sequential Monte Carlo methods in practice*. Springer Verlag, 2001.
- [13] S. Thrun, "Probabilistic robotics," *Communications of the ACM*, vol. 45, no. 3, pp. 52–57, 2002.
- [14] D. Fox, W. Burgard, F. Dellaert, and S. Thrun, "Monte carlo localization: Efficient position estimation for mobile robots," in *Proceedings of the National Conference on Artificial Intelligence*. John Wiley & sons Ltd, 1999, pp. 343–349.
- [15] C. Kwok, D. Fox, and M. Meila, "Adaptive real-time particle filters for robot localization," in *Robotics and Automation, 2003. Proceedings. ICRA'03. IEEE International Conference on*, vol. 2. IEEE, 2003, pp. 2836–2841.
- [16] M. Montemerlo, S. Thrun, D. Koller, and B. Wegbreit, "Fastslam 2.0: An improved particle filtering algorithm for simultaneous localization and mapping that provably converges," in *International Joint Conference on Artificial Intelligence*, vol. 18, 2003, pp. 1151–1156.
- [17] R. Sim, P. Elinas, M. Griffin, and J. Little, "Vision-based slam using the rao-blackwellised particle filter," in *IJCAI Workshop on Reasoning with Uncertainty in Robotics*, 2005, pp. 9–16.
- [18] A. Howard, "Multi-robot simultaneous localization and mapping using particle filters," *The International Journal of Robotics Research*, vol. 25, no. 12, pp. 1243–1256, 2006.
- [19] A. Notsu, K. Honda, and H. Ichihashi, "Proposed particle-filtering method for reinforcement learning," in *Fuzzy Systems (FUZZ), 2011 IEEE International Conference on*. IEEE, 2011, pp. 1755–1718.
- [20] A. Notsu, K. Honda, H. Ichihashi, Y. Komori, and Y. Iwamoto, "Improvement of particle filter for reinforcement learning," *International Conference on Machine Learning and Applications*, vol. 1, pp. 454–457, 2011.
- [21] K. Samejima, K. Doya, Y. Ueda, and M. Kumura, "Estimating internal variables and parameters of a learning agent by a particle filter," in *Advances in Neural Information Processing Systems*, vol. 16. NIPS, 2004, pp. 1335–1342.
- [22] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer New York, 2006, (pages 644–646).
- [23] N. Gordon, D. Salmond, and A. Smith, "Novel approach to nonlinear/non-gaussian bayesian state estimation," in *Radar and Signal Processing, IEE Proceedings F*, vol. 140, no. 2. IET, 1993, pp. 107–113.