

Lifelong Learning Applications to Mobile Robotics

David Isele
Univ. of Pennsylvania
isele@seas.upenn.edu

José Marcio Luna
Univ. of Pennsylvania
joseluna@seas.upenn.edu

Eric Eaton
Univ. of Pennsylvania
eeaton@seas.upenn.edu

Gabriel de la Cruz
Washington State Univ.
gabriel.delacruz@wsu.edu

James Irwin
Washington State Univ.
james.irwin@wsu.edu

Brandon Kallaher
Washington State Univ.
brandon.kallaher@wsu.edu

Matthew E. Taylor
Washington State Univ.
taylor@eecs.wsu.edu

Marcio's comment: We need to decide a good title.

ABSTRACT

Learning controllers for multiple systems is often an expensive process when controllers for each system are learned individually. Advances in lifelong learning suggest that information between systems can be shared, improving the quality of the controllers that are learned. However these results have been largely theoretical, with applications limited to benchmark problems with known dynamics. We show that these methods can be extended to robotic platforms. Particularly we validate our assumptions for transfer learning between tasks with unknown dynamics in order to carry out a disturbance rejection problem. We view this as early work leading up to learning robust fault tolerant control.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning; I.2.9 [Artificial Intelligence]: Robotics

General Terms

algorithms, experimentation

Keywords

lifelong learning, policy gradients, reinforcement learning, robotics

1. INTRODUCTION

In control systems, a perfect model of the dynamics of the system is often necessary to guarantee the stabilization of the system. This can be problematic in complicated systems where the dynamics are difficult to model or require information that is not available to the designer.

Policy search approaches have been proposed to deal with the design of controllers in model-free applications. However, it is difficult to make claims of robustness or generalization of the learned controllers **Marcio's Comment: We need citations for this.**


Our goal is to learn fault tolerant control in multi-agent systems. In order to approach this problem we begin by

Appears in: *Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2016)*, John Thangarajah, Karl Tuyls, Stacy Marsella, Catholijn Jonker (eds.), May 9–13, 2016, Singapore. Copyright © 2016, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

developing a method that can accomplish disturbance rejection across multiple robots. Given a collection of robots each with a different unknown disturbance, learning a unique control policy for each robot in the system can be very costly. One approach to reduce the amount of learning is to share information between robots.

Lifelong learning [10] is a promising approach for accomplishing information sharing between different robots. It works on-line, allowing the different systems to be encountered consecutively so a priori knowledge of all the different systems is not required at the start of training. Also it preserves and possibly improves the models encountered early on, in contrast to transfer methods which only optimize performance on the target system.

However most recent work in lifelong learning [10, 2, 3] has been theoretical, using simulations of benchmark problems with known dynamics to demonstrate knowledge sharing. The contribution of this work is to present the first results of lifelong learning on robotic systems. We do this by applying the PG-ELLA framework [2] to the problem of disturbance rejection.

[ **Marcio:** The introduction needs more work, please don't forget to clarify our main contributions/goals and please provide the organization of the paper.]

2. RELATED WORK

Building mathematical models that describe the behavior of physical systems is common practice to analyze, predict and control their behavior to fulfill specific goals. Among the well known techniques for modeling physical systems we find partial, ordinary differential and difference equations [?, ?], and Discrete Event Systems (DES) such as queueing and Petri networks [?, ?, ?]. Control systems theory uses differential and difference equations to model the mechanics of physical systems. In control systems a validated model provides ways of assessing the stability and stabilizability of the system to be able to control its behavior.

Some of the typical problems in control systems are regulation, trajectory tracking, disturbance rejection and robustness among many others [?, ?, ?]. All these problems are associated with the analysis of the stabilizability of the system, as well as the design of controllers to stabilize it. These controllers consist of theoretical artifacts that would allow the solution of control problems such as the aforementioned ones.

The disturbance rejection problem consists of implementing a controller that allows the plant to fulfill the desired task while compensating for a disturbance that modifies its nom-

inal dynamics. As long as there is an accurate model of the plant, several mathematical artifacts have been provided to handle constant, constant and unknown, time-varying and even stochastic disturbances [?, ?, ?]. However, things get more complicated when no model is provided, even for simple disturbances. In a model-free setting, the goal is to design a controller that stabilizes a system whose model is not available due to complex internal iterations, uncertainty in the system, event-based dynamics and technological limitations.

[**☛ Marcio:** From here on this section needs more specifics about model-free approaches. Please make sure you cite relevant and recent literature, e.g., the Pilco paper.]

Reinforcement learning [5] is often utilized to learn controllers in a model-free settings. Amongst reinforcement learning algorithms, policy gradient (PG) methods [11] are popular in robotic applications since they accommodate continuous state/action spaces and can scale well to high dimensional spaces.

It has been shown that PG methods can be used in a lifelong learning setting [2, 3], however these works focused largely on theory, using benchmark simulations to demonstrate their results. While there are examples of lifelong learning on robots, they tend to focus in skill refinement on a single robot [4, 12] rather than sharing information across multiple robots as we do in our work.

3. BACKGROUND

Our approach works by sharing knowledge between different robotic systems. The policy for each system is learned by reinforcement learning. In this section we cover the mathematical framework that supports our experiments on lifelong learning.

3.1 Reinforcement Learning

A reinforcement learning (RL) agent must select sequential actions to maximize its expected return. RL approaches do not require previous knowledge of the system dynamics, instead, the control policies are learned through the interactions with the system. RL problems are typically formalized as Markov Decision Processes (MDPs) with the form $\langle \mathcal{X}, \mathcal{A}, P, R, \gamma \rangle$ where $\mathcal{X} \subset \mathbb{R}^{d_x}$ is the set of states, \mathcal{A} is the set of actions, $P : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$ is the state transition probability describing the systems dynamics, $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function and $\gamma \in [0, 1]$ is the reward discount factor. At each time step h , the agent is in the state $\mathbf{x}_h \in \mathcal{X}$ and must choose an action $\mathbf{a}_h \in \mathcal{A}$ so that it transitions to a new state \mathbf{x}_{h+1} with state transition probability $P(\mathbf{x}_{h+1}|\mathbf{x}_h, \mathbf{a}_h)$, yielding a reward r_h according to R . The action is selected according to a policy $\pi : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ which specifies a probability distribution over actions given the current state. The goal of RL is to find the optimal policy π^* that maximizes the expected reward.

We use a class of RL algorithms known as Policy Gradient (PG) methods [11], which are particularly well suited for solving high dimensional problems with continuous state and action spaces, such as robotic control [9].

The goal of PG is to use gradient steps to optimize the expected average return:

$$\mathcal{J}(\boldsymbol{\theta}) = \int_{\mathbb{T}} p_{\boldsymbol{\theta}}(\tau) \mathcal{R}(\tau) d\tau, \quad (1)$$

where \mathbb{T} is the set of all trajectories and $\mathcal{R}(\tau)$ is the average

per-step reward, specifically:

$$p_{\boldsymbol{\theta}} = \prod_{h=0}^H p(\mathbf{x}_{h+1}|\mathbf{x}_h, \mathbf{a}_h) \pi(\mathbf{a}_h, \mathbf{x}_h),$$

$$\mathcal{R}(\tau) = \frac{1}{H} \sum_{h=0}^H r(\mathbf{s}_h, \mathbf{a}_h, \mathbf{s}_{h+1}).$$

Most PG methods (e.g. episodic REINFORCE [13], PoWER [6], and Natural Actor Critic [9]) optimize the policy by maximizing a lower bound on the return, comparing trajectories generated by different candidate policies π . In this particular application, the PG method we use in our experiments is finite differences (FD) [?] which optimizes the return directly.

3.2 Finite Differences for Policy Search

The local optimization around an existing policy π parameterized by a parameter matrix $\boldsymbol{\theta}$ is carried out by computing changes in the policy parameters $\Delta\boldsymbol{\theta}$ that will increase the expected reward, thus producing the iterative update,

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + \Delta\boldsymbol{\theta}_m.$$

Gradient-based methods for policy updates follow the gradient of the expected return \mathcal{J} given a step-size δ ,

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + \delta \nabla_{\boldsymbol{\theta}} \mathcal{J}.$$

In FD gradients, we have a set of n perturbed policy parameters which are used to estimate the effect of a change in policy parameters:

$$\Delta\hat{\mathcal{J}}_{\mathbf{p}} \approx \mathcal{J}(\boldsymbol{\theta}_m + \Delta\boldsymbol{\theta}_{\mathbf{p}}) - \mathcal{J}_{ref},$$

where $\Delta\boldsymbol{\theta}_{\mathbf{p}}$ are the individual perturbations for $\mathbf{p} = [1, \dots, n]$, $\Delta\hat{\mathcal{J}}_{\mathbf{p}}$ is the estimate of their effect on the return, and the \mathcal{J}_{ref} is a reference return which is usually taken as the return of the unperturbed parameters. By using linear regression we get an approximation of the gradient,

$$\nabla_{\boldsymbol{\theta}} \mathcal{J} \approx \left(\Delta\boldsymbol{\theta}^T \Delta\boldsymbol{\theta} \right)^{-1} \Delta\boldsymbol{\theta}^T \Delta\hat{\mathcal{J}}_{\mathbf{p}},$$

where $\Delta\hat{\mathcal{J}}_{\mathbf{p}}$ contains all the stacked samples of $\Delta\hat{\mathcal{J}}_{\mathbf{p}}$ and $\Delta\boldsymbol{\theta}$ contains the stacked perturbations $\Delta\boldsymbol{\theta}_{\mathbf{p}}$. This approach is sensitive to the type and magnitude of the perturbations, as well as to the step size δ . Since the number of perturbations needs to be as large as the number of parameters, this method is considered to be noisy and inefficient for problems with large sets of parameters.

In our experiments, the policy is represented as a function defined over a parameter matrix $\boldsymbol{\theta} \in \mathbb{R}^{d_{\boldsymbol{\theta}} \times M}$, that provides the gains of the M control inputs of the system. Our goal is to optimize the expected average return of Equation 1.

In order to share information across different learned policies, we incorporate the PG learning process into a lifelong machine learning setting.

3.3 Lifelong Machine Learning

Lifelong learning focuses on learning a set of tasks consecutively while performing well across all tasks. Given a round $t = 1, \dots, t_{\max}$ a task $T^{(t)}$ is observed. [**☛ Eric:** Please don't change notation between papers unless absolutely necessary. Let's use the exact same notation as in the recent IJCAI paper.] In our setting, each task corresponds to a reinforcement learning problem for an individual robot.

We assume that the model associated to $T^{(t)}$ is parameterized by a parameter $\theta^{(t)} \in \mathbb{R}^{d_{\theta} \times M}$. The ideal goal is that prior knowledge about tasks $T^{(1)}, \dots, T^{(t-1)}$ should provide enough information so that the lifelong learning algorithm performs better and faster on $T^{(t)}$ while being able to scale as the number of tasks increases.

Following work in both multi-task [?] and lifelong learning [10], we assume there is a shared basis $\mathbf{L} \in \mathbb{R}^{d_x \times l}$ and a sparse weight vector $\mathbf{s}^{(t)} \in \mathbb{R}$, so that the policy parameters $\theta^{(t)}$ are given by,

$$\theta^{(t)} = \mathbf{L} \mathbf{s}^{(t)}.$$

Using the return function in (1) we propose the following multi-task objective function,

$$\arg \min_{\mathbf{L}, \mathbf{s}} \frac{1}{|T|} \sum_t \left[-\mathcal{J}(\theta^{(t)}) + \lambda \|\mathbf{s}^{(t)}\|_1 \right] + \mu \|\mathbf{L}\|_F^2,$$

where \mathbf{S} is the set of the sparse vectors $\mathbf{s}^{(t)}$, the L1 norm of $\|\mathbf{s}^{(t)}\|_1$ provides sparse code for $\mathbf{s}^{(t)}$ and the Frobenious norm in $\|\mathbf{L}\|_F^2$ provides regularization of \mathbf{L} . The coefficients μ and $\lambda \in \mathbb{R}$ are weights for the regularization and sparsity respectively. The learning objective function is approximated by a second order Taylor expansion around an estimate $\alpha^{(t)}$ of the single task policy parameters of task t . The optimization problem is solved by using the on-line ELLA algorithm introduced in [10] and extended to reinforcement learning in [2]. The optimization problem is solved by incrementally updating \mathbf{L} by the following the update equations,

$$\begin{aligned} \mathbf{s}^{(t)} &\leftarrow \arg \min_{\mathbf{s}} \left\| \alpha^{(t)} - \mathbf{L} \mathbf{s}^{(t)} \right\|_{\Gamma^{(t)}}^2 + \mu \|\mathbf{s}\|_1, \\ \mathbf{A} &\leftarrow \mathbf{A} + \left(\mathbf{s}^{(t)} \mathbf{s}^{(t)\top} \right) \otimes \Gamma^{(t)}, \\ \mathbf{b} &\leftarrow \mathbf{b} + \text{vec} \left(\mathbf{s}^{(t)} \otimes \left(\theta^{(t)\top} \Gamma^{(t)} \right) \right), \\ \mathbf{L} &\leftarrow \text{mat} \left(\left(\frac{1}{T} \mathbf{A} + \lambda \mathbf{I}_{l \times d_{\theta}, l \times d_{\theta}} \right)^{-1} \frac{1}{T} \mathbf{b} \right). \end{aligned} \quad (2)$$

where $\|\mathbf{v}\|_{\mathbf{A}}^2 = \mathbf{v}^\top \mathbf{A} \mathbf{v}$ and $\Gamma^{(t)}$ is the Hessian of the PG objective function, and \mathbf{A} and \mathbf{b} are initialized to zero matrices.

4. EXPERIMENTS

We evaluated our approach by modeling the control policies for different Turtlebot systems [?]. Turtlebots are an affordable robotic platform that uses the Robotic Operating System (ROS) in its hydro version, which is compatible with the Gazebo simulator [?, ?]. We artificially induce a random and constant disturbance to the control signal to emulate a bias on the angular velocity of each robot. Each turtlebot has a different constant disturbance drawn uniformly from $[-0.1, 0.1] \subset \mathbb{R}$ and measured in m/s .

The Gazebo simulator considers the kinematics and mechanics of the system. In a first attempt to explore disturbance rejection as a preamble of fault tolerant control applications, the Turtlebot should learn how to drive itself from an initial to a goal point. It should accomplish this simple task while being affected by a constant and unknown angular disturbance. Thus the robot is enforced to compensate for the induced failure. It is worth mentioning that the difficulty of the disturbance will be increased by assuming time varying disturbance and stochastic disturbances later on.

We assume we have little knowledge of the Turtlebot model, so we use the oversimplified kinematic model provided in [1].

Algorithm 1 PG-ELLA (k, λ, μ)

```

1:  $T \leftarrow 0$ 
2:  $\mathbf{A} \leftarrow \text{zeros}_{k \times d, k \times d}$ ,  $\mathbf{b} \leftarrow \text{zeros}_{k \times d, 1}$ 
3:  $\mathbf{L} \leftarrow \text{RandomMatrix}_{d, k}$ ,  $\mathbf{D} \leftarrow \text{RandomMatrix}_{m, k}$ 
4: while some task  $t$  is available do
5:   if isNewTask( $t$ ) then
6:      $T \leftarrow T + 1$ 
7:      $(\mathbb{T}^{(t)}, R^{(t)}) \leftarrow \text{getRandomTrajectories}()$ 
8:   else
9:      $(\mathbb{T}^{(t)}, R^{(t)}) \leftarrow \text{getTrajectories}(\alpha^{(t)})$ 
10:     $\mathbf{A} \leftarrow \mathbf{A} - (\mathbf{s}^{(t)} \mathbf{s}^{(t)\top}) \otimes \Gamma^{(t)}$ 
11:     $\mathbf{b} \leftarrow \mathbf{b} - \text{vec}(\mathbf{s}^{(t)\top} \otimes (\alpha^{(t)\top} \Gamma^{(t)}))$ 
12:   end if
13:   Compute  $\alpha^{(t)}$  and  $\Gamma^{(t)}$  from  $\mathbb{T}^{(t)}$ 
14:    $\mathbf{s}^{(t)} \leftarrow \arg \min_{\mathbf{s}} \left\| \alpha^{(t)} - \mathbf{L} \mathbf{s}^{(t)} \right\|_{\Gamma^{(t)}}^2 + \mu \|\mathbf{s}\|_1$ 
15:    $\mathbf{A} \leftarrow \mathbf{A} + (\mathbf{s}^{(t)} \mathbf{s}^{(t)\top}) \otimes \Gamma^{(t)}$ 
16:    $\mathbf{b} \leftarrow \mathbf{b} + \text{vec}(\mathbf{s}^{(t)\top} \otimes (\alpha^{(t)\top} \Gamma^{(t)}))$ 
17:    $\mathbf{L} \leftarrow \text{mat} \left( \left( \frac{1}{T} \mathbf{A} + \lambda \mathbf{I}_{k \times d, k \times d} \right)^{-1} \frac{1}{T} \mathbf{b} \right)$ 
18:   for  $t \in \{1, \dots, T\}$  do:  $\theta^{(t)} \leftarrow \mathbf{L} \mathbf{s}^{(t)}$ 
19: end while
```

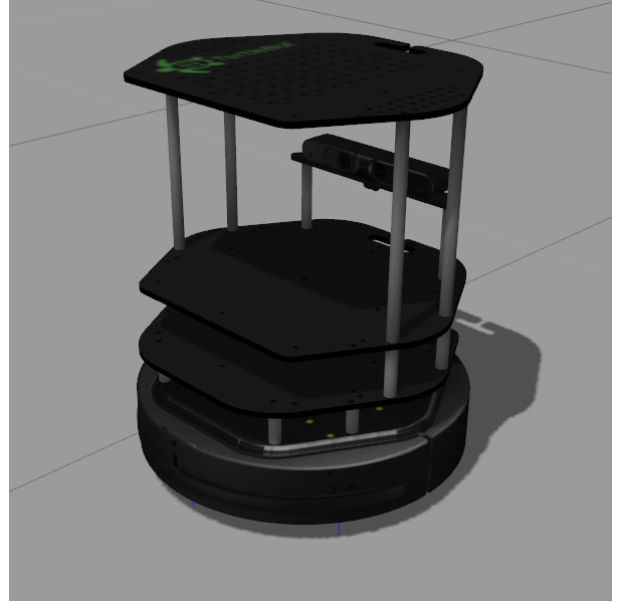


Figure 1: Turtlebot model in Gazebo Simulator.

Marcio's comment: The reason why figures are usually squared in papers is because authors prefer to fill the largest space possible so they look as big as possible. I suggest to try to enlarge the figures. Their purpose is to be illustrative and that sometimes implies size.

In the kinematic model, the estate space is given by $\mathbf{x} \in \mathbb{R}^3$ and the action space is described by $\mathbf{a} \in \mathbb{R}^2$. Notice that the model in [1] just considers the kinematics of a unicycle in polar coordinates so we do not consider the dynamics of the system, *i.e.*, we are neglecting model parameters such

as mass, damping and friction coefficients, as well as inputs such as forces and torques. Then, the nonlinear policy is derived neglecting the dynamics and just assuming simple kinematics, therefore, taking into account that our action is given by $\alpha = \theta^T \phi(x)$ where $x = (\rho, \gamma, \psi)$ as illustrated in Fig. 3 we propose the following nonlinear gain vector structure, The names of the state variables in the figure don't correspond to the names used in the control gain, namely, α, ψ, κ . [**Eric:** Remember, α is the estimated single-task policy parameters, so don't overload the notation.]

$$\phi(x) = \begin{pmatrix} \rho \cos(\gamma) \\ \frac{\gamma}{\cos(\gamma) \sin(\gamma)} (\gamma + \psi) \\ 1 \end{pmatrix} \quad (3)$$

where ρ, γ and ψ are indicated in Fig. 2.

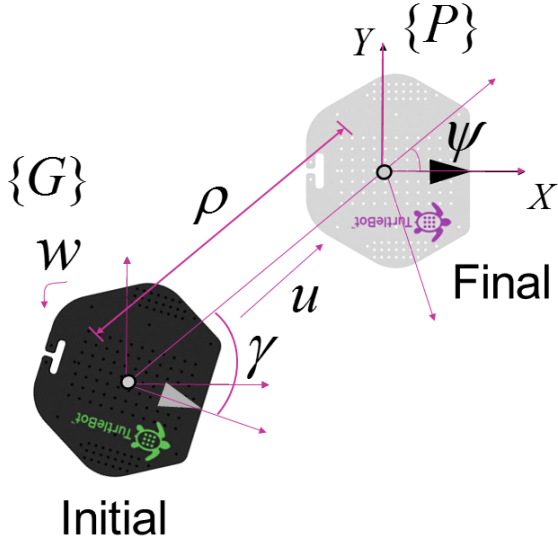


Figure 2: State variables of simplified go-to-goal problem. **Marcio's comments:** I usually prefer .eps images, they look way better and you can amplify them and they have really good resolution. I know they're a pain in the back, but is worthy. However, this is your paper so it's your call.

The state is a non-linear transformation of the position and heading angle. This implementation was entirely done in Python.

Although this is true, these kind of assertions are not appreciated unless there's a good reason for it we could provide. In that case is better not to mention the fact, unless we have a hypothesis or a conclusion of what's going on. The reason why this is incorrect has to do with the fact that the original authors of REINFORCE and NAC would desire an explanation about why we are questioning the correctness of their algorithm, given that a lot of papers rely on their result, therefore, we look as if we don't know what we're doing, and that isn't the case. We just need to somehow come up with a formal proof that our result is better. We were not able to learn a policy on this non-linear state transformation using either the natural actor critic [9] or episodic

REINFORCE [13] policy gradient methods. [**Eric:** WHY are we using FD? I.e., what are its advantages in this application over REINFORCE and NAC?]

4.1 Methodology

Marcio's comments: This is just a suggestion. I usually divide the experimental section in two subsections, namely, 'experimental design' and 'Results'. In the first one, I provide the details of the experiment such as parameters, programming language, hardware specs, etc. In the second one I present the result and the analysis of the results. I believe this section still lacks more analysis. We start by generating 20 robots, each with a different constant disturbance and a unique goal, both selected uniformly. The robots are run for 20 learning steps of FD, each learning step consists of 15 roll-outs of 50 time steps each. Note that 20 learning steps is fewer iterations than is required for any system to converge to a good controller. The policy that is learned after 20 iterations of FD is used as θ^* for PG-ELLA. We use a learning rate of 1×10^{-6} .

We learn our PG-ELLA knowledge repository \mathcal{L} and sparse representation $s^{(t)}$ using the update equations given by (2). Tasks are encountered randomly with repetition and learning stops once every task has been seen once. We select $k = 8$ to be the number of columns, and use sparsity coefficient $\mu = 1 \times 10^{-3}$ and regularization coefficient $\lambda = 1 \times 10^{-8}$.

We then compare the policy reconstructed from PG-ELLA against the policy that was learned after 20 iterations of FD by comparing the learning curves that result from running FD for an additional 80 learning iterations. Performance is averaged over 6 trials for all 20 robots. In Fig. 3 we see that PG-ELLA is successfully able to reconstruct the policies and provide an additional benefit of positive transfer. Note that these are preliminary results and we suspect further refinements will enable us to achieve larger amounts of transfer.

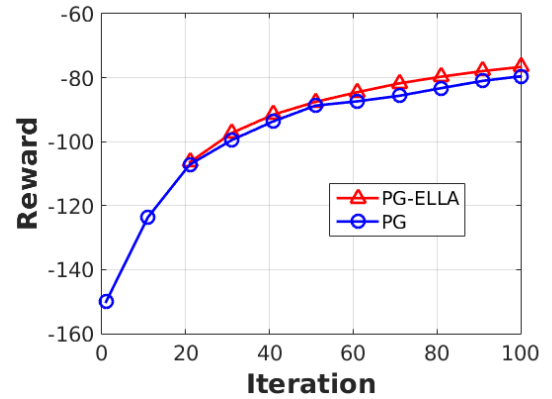


Figure 3: Learning curves for FD and PG-ELLA using FD to transfer information between tasks.

5. CONCLUSIONS

Marcio's Comment: David please finish the conclusions We demonstrate the use of lifelong learning for disturbance rejection on Turtlebots. This is intended to lay the foundation for creating fault tolerant control on multi-agent systems.

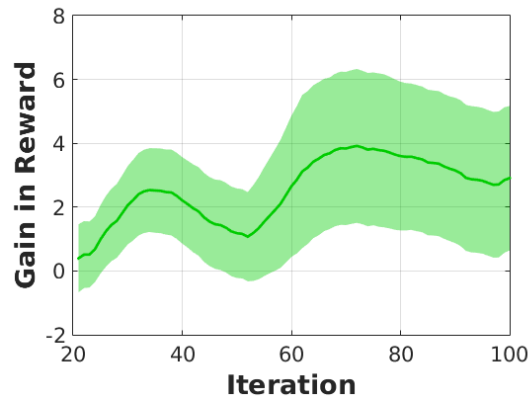


Figure 4: The positive transfer achieved by using ELLA.

Acknowledgments

the grant goes here

6. ADDITIONAL AUTHORS

REFERENCES

- [1] M. Aicardi, G. Casalino, A. Balestrino, and A. Bicchi. Closed loop smooth steering of unicycle-like vehicles. In *Decision and Control, 1994., Proceedings of the 33rd IEEE Conference on*, volume 3, pages 2455–2458. IEEE, 1994.
- [2] H. Bou Ammar, E. Eaton, and P. Ruvolo. Online multi-task learning for policy gradient methods. *International Conference on Machine Learning*, 2014.
- [3] H. Bou Ammar, E. Eaton, P. Ruvolo, and M. E. Taylor. Unsupervised cross-domain transfer in policy gradient reinforcement learning via manifold alignment. *International Joint Conference on Artificial Intelligence*, 2015.
- [4] A. Kleiner, M. Dietl, and B. Nebel. Towards a life-long learning soccer agent. In *RoboCup 2002: Robot Soccer World Cup VI*, pages 126–134. Springer, 2002.
- [5] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, page 0278364913495721, 2013.
- [6] J. Kober and J. Peters. Policy search for motor primitives in robotics. *Advances in Neural Information Processing Systems*, pages 849–856, 2009.
- [7] A. Kumar and H. Daumé. Learning task grouping and overlap in multi-task learning. *International Conference on Machine Learning*, pages 1383–1390, 2012.
- [8] J. Peters and S. Schaal. Policy gradient methods for robotics. pages 2219–2225, 2006.
- [9] J. Peters and S. Schaal. Natural actor-critic. *Neurocomputing*, 71(7):1180–1190, 2008.
- [10] P. Ruvolo and E. Eaton. ELLA: An efficient lifelong learning algorithm. *International Conference on Machine Learning*, 28:507–515, 2013.
- [11] R. S. Sutton, D. A. McAllester, S. P. Singh, and

Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems*, 99:1057–1063, 1999.

- [12] S. Thrun and T. M. Mitchell. *Lifelong robot learning*. Springer, 1995.

- [13] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.