

Scalable Lifelong Reinforcement Learning

Yusen Zhan

*The School of Electrical Engineering and Computer Science, Washington State University,
Pullman, WA 99163, USA*

Haitham Bou Ammar

Prowler i.o., Cambridge, United Kingdom

Matthew E. Taylor*

*The School of Electrical Engineering and Computer Science, Washington State University,
Pullman, WA 99163, USA*

Abstract

Lifelong reinforcement learning provides a successful framework for agents to learn multiple consecutive tasks sequentially. Current methods, however, suffer from scalability issues when the agent has to solve a large number of tasks. In this paper, we remedy the above drawbacks and propose a novel scalable technique for lifelong reinforcement learning. We derive an algorithm which assumes the availability of multiple processing units and computes shared repositories and local policies using only local information exchange. We then show an improvement to reach a *linear convergence rate* compared to current lifelong policy search methods. Finally, we evaluate our technique on a set of benchmark dynamical systems and demonstrate learning speed-ups and reduced running times.

Keywords: Reinforcement Learning, Lifelong Learning, Distributed Optimization, Transfer Learning

*Corresponding author

Email addresses: yusen.zhan@wsu.edu (Yusen Zhan), haitham@prowler.io (Haitham Bou Ammar), taylor@eeecs.wsu.edu (Matthew E. Taylor)

1. Introduction

Reinforcement learning (RL) provides the ability to solve sequential decision-making problems with limited feedback. Applications with these characteristics range from robotics control [1] to personalized medicine [2, 3]. Though successful, typical RL methods require a substantial amount of experience before
5 acquiring acceptable behavior. The cost of obtaining such experience, however, can be prohibitively expensive in terms of time and data [4]. Also, these costs only worsen when considering multiple tasks¹. Transfer learning² and multi-task learning [7, 8] have been developed to remedy these problems by allowing agents
10 to reuse knowledge from other tasks. Unfortunately, both these techniques suffer from scalability problems as the number of tasks considered grows large. In a recent attempt to target these issues, Bou Ammar *et al.* proposed PG-ELLA, an online hybrid of the above two paradigms—transfer and multi-task learning [9]. PG-ELLA decomposes a task’s policy parameters into a shared latent
15 repository \mathbf{L} and task specific coefficients \mathbf{s}_t , one for each task. Consequently, it allows for knowledge transfer between tasks (using the latent repository \mathbf{L}) while scaling multi-task learning by streaming problems online. In the original paper, the authors show a closed form solution to the shared repository.

In a lifelong setting, the agent deals with a sequence of tasks, potentially leading to computational intractability when the agent observes and learns many
20 tasks (e.g., in big data scenario, the agent may observe many tasks in its life). Unfortunately, PG-ELLA suffers when considering a large number of tasks or dimensions — determining the shared knowledge-base \mathbf{L} becomes intractable — due to two inefficiencies. The first inefficiency is that computing the expansion’s
25 operating point amounts to solving a local RL problem described by the current task’s observed trajectories. The second inefficiency reducing PG-ELLA’s

¹The multiple tasks learning tried to learn a single policy or strategy to deal with a set of tasks. See [5] for details.

²The idea of transfer learning is to reuse previous trained information to speed up learning process. See [6] for details.

scalability arises when updating the shared repository \mathbf{L} .

In this paper, we remedy the above drawbacks of lifelong policy search and propose a scalable method capable of handling large number of tasks and high-dimensional policies, while ensuring meaningful parameters at each iteration of the algorithm (i.e., each step conveys more information than a single gradient step, as done in PG-ELLA). Our method assumes multiple processing units and derives a novel algorithm for lifelong policy search. Such an assumption leads us to better convergence rates and more accurate local policies. Using augmented Lagrangian methods, we propose two algorithms for computing the local operating point and the shared repository, while allowing for scalable solutions. We show that such updates can be computed locally by each processing unit in closed form under broad assumptions. We theoretically and empirically analyze the performance of these new techniques. On the theoretical side, we demonstrate superiority to PG-ELLA by proving *dimensionality-free linear convergence* in both determining local policies and updating the shared repository. On the empirical side, we show our algorithm outperforms existing algorithms on a set of 5 domains with 50 variations, including the control of a simulated helicopter. In a final set of experiments, we also demonstrate the improved generalization capabilities of this new method on unobserved tasks and report a decrease in learning time, relative to current techniques. The contributions of this paper can be summarized as: *i)* deriving a novel scalable algorithm for lifelong policy search, *ii)* acquiring linear convergence rate in both steps of the method, *iii)* demonstrating the effectiveness of our technique on 5 dynamical systems with 50 variations each, and *iv)* demonstrating learning speed-ups on unobserved tasks.

The rest of the paper is organized as follows: Section 2 introduces the background; Section 3 provides some related works; In section 4, we define our lifelong policy search problem; Section 5 proposes the scalable lifelong policy search method; In section 6, it shows the experimental results and section 7 summarizes the whole paper.

2. Reinforcement Learning

In reinforcement learning (RL) an agent must sequentially select actions to maximize its total expected pay-off. These problems are typically formalized as Markov decision processes (MDPs) $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma \rangle$, where $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{U} \subseteq \mathbb{R}^m$ denote the state and action spaces. $\mathcal{P} : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \rightarrow [0, 1]$ represents the transition probability governing the dynamics of the system, $\mathcal{R} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ is the reward function quantifying the performance of the agent and $\gamma \in [0, 1)$ is a discount factor specifying the degree to which rewards are discounted over time. At each time step m , the agent is in state $\mathbf{x}_m \in \mathcal{X}$ and must choose an action $\mathbf{u}_t \in \mathcal{U}$, transitioning it to a successor state $\mathbf{x}_{m+1} \sim p(\mathbf{x}_{m+1} | \mathbf{x}_m, \mathbf{u}_m)$ as given by \mathcal{P} and yielding a reward $r_{m+1} = \mathcal{R}(\mathbf{x}_m, \mathbf{u}_m)$. A policy $\pi : \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ is defined as a probability distribution over state-action pairs, where $\pi(\mathbf{u}_m | \mathbf{x}_m)$ denotes the probability of choosing action \mathbf{u}_t in state \mathbf{x}_m .

Policy gradients [1, 4] are a class of reinforcement learning algorithms that have shown successes in solving complex robotic problems [1]. Such methods represent the policy $\pi_{\theta}(\mathbf{u}_m | \mathbf{x}_m)$ by an unknown vector of parameters $\theta \in \mathbb{R}^d$. The goal is to determine the optimal parameters θ^* that maximize the expected average pay-off:

$$\mathcal{J}(\theta) = \int_{\tau} p_{\theta}(\tau) \mathfrak{R}(\tau) d\tau,$$

where $\tau = [\mathbf{x}_{0:T}, \mathbf{u}_{0:T}]$ denotes a trajectory over a possibly infinite horizon T . The probability of acquiring a trajectory, $p_{\theta}(\tau)$, under the policy parameterization $\pi_{\theta}(\cdot)$ and average per-time-step return $\mathfrak{R}(\tau)$ is given by:

$$p_{\theta}(\tau) = p_0(\mathbf{x}_0) \prod_{m=0}^{M-1} p(\mathbf{x}_{m+1} | \mathbf{x}_m, \mathbf{u}_m) \pi_{\theta}(\mathbf{u}_m | \mathbf{x}_m), \quad \mathfrak{R}(\tau) = \frac{1}{M} \sum_{m=0}^M r_{m+1},$$

with an initial state distribution $p_0 : \mathcal{X} \rightarrow [0, 1]$. Policy gradient methods, such as episodic REINFORCE [10] and Natural Actor Critic [11, 12], typically employ a lower-bound on the expected return $\mathcal{J}(\theta)$ for fitting the unknown policy parameters θ . To achieve this, such algorithms generate trajectories using the current policy θ , and then compares performance with a new parameterization

$\tilde{\theta}$. As detailed in [1], the lower bound on the expected return can be attained using Jensen’s inequality and the concavity of the logarithm:

$$\begin{aligned} \log \mathcal{J}(\tilde{\theta}) &= \log \int_{\tau} p_{\tilde{\theta}}(\tau) \mathfrak{R}(\tau) d\tau \geq \int_{\tau} p_{\theta}(\tau) \mathfrak{R}(\tau) \log \frac{p_{\tilde{\theta}}(\tau)}{p_{\theta}(\tau)} d\tau + \text{cnst.} \\ &\propto -\text{KL}(p_{\theta}(\tau) \mathfrak{R}(\tau) \| p_{\tilde{\theta}}(\tau)) = \mathcal{J}_{\mathcal{L}, \theta}(\tilde{\theta}), \end{aligned}$$

70 where $\text{KL}(p(\tau) \| q(\tau)) = \int_{\tau} p(\tau) \log \frac{p(\tau)}{q(\tau)} d\tau$. Thus, we can directly optimize the lower bound instead of the original objective.

3. Related Work

Transfer learning aims to improve the learning speed of an agent on a new target task by transferring knowledge learned from one or more previous source
75 tasks [6]. However, most current transfer learning methods only consider the target tasks instead of all previously visited tasks. In contrast, multi-task learning (MTL) methods optimize performance over all tasks by training models over all possible tasks [7, 8], but are computationally expensive in lifelong learning setting since the agent must learn a large number of tasks over time. [13, 14].

80 There are many works that consider parallelizing reinforcement learning. For example, Caarls and Schuitema apply parallel on-line temporal Temporal Difference (TD) learning to motor control [15]. Shixiang *et al.* proposed an asynchronous parallel method to deep reinforcement learning [16]. Yahya *et al.* introduce a distributed asynchronous guided policy search method to deep
85 reinforcement learning [17]. Sergey *et al.* also employed Bregman ADMM, which is a centralized ADMM method, to guided policy search setting [18]. However, all of the above methods are based on standard reinforcement learning instead of lifelong learning setting.

PG-ELLA is a recent lifelong policy gradient RL algorithm that can effi-
90 ciently learn multiple tasks consecutively while sharing knowledge between task policies to accelerate learning [9]. In fact, PG-ELLA, which is a centralized method, can be viewed as one agent or node case of the algorithm we propose in this paper. MTL for policy gradients has also been explored by Deisenroth

et al. [19] through customizing a single parameterized controller to individual
 95 tasks that differ only in the reward function. Another closely related work is
 on hierarchical Bayesian MTL [20], which can learn RL tasks consecutively, but
 unlike our approach, requires discrete states and actions. Snel and Whiteson’s
 representation learning approach is also related, but limited to a potential func-
 tion representation [21]. GO-MTL is also a multi-task learning algorithm that
 100 focuses on multi-task supervised learning [22]. All of these MTL methods suffer
 from scalability issues due to their dependence on centralized methods.

4. Lifelong Policy Search

We adopt the lifelong learning framework previously introduced in [9]. An agent has to master multiple RL tasks while supporting transfer to improve learning. Let \mathcal{T} denote the set of tasks in which each element is an MDP. At any time, the learner may face any of the previously seen tasks — it must maximize its performance across all tasks. The goal is to learn optimal policies, $\pi_{\theta_1^*}, \dots, \pi_{\theta_{|\mathcal{T}|}^*}$ for all tasks with corresponding parameters $\theta_1^*, \dots, \theta_{|\mathcal{T}|}^*$, where $\theta_j^* \in \mathbb{R}^d$ parameterizes the policy for task j . For each task j , the learner samples a set of n_j trajectories $\mathbb{T}^{(j)} = \{\tau_j^{(1)}, \dots, \tau_j^{(n_j)}\}$ from the environment by predefined policies, and each trajectory has a length $M^{(j)}$. To transfer knowledge between tasks, we adopt the shared knowledge-base representation, where the policy parameters for each task are represented as a linear combination of a shared latent basis $\mathbf{L} \in \mathbb{R}^{d \times p}$ and coefficient vectors $\mathbf{s}_j \in \mathbb{R}^p$, i.e., $\theta_j = \mathbf{L}\mathbf{s}_j$. Here, each column of \mathbf{L} encodes a set of transferable knowledge that is tailored to suit the peculiarities of each task using \mathbf{s}_j . Such a decomposition has been widely used in the multi-task learning literature [9, 13, 22, 23]. Hence, the optimization problem for a total of t tasks can be written as:

$$\min_{\mathbf{L}, \mathbf{s}_1: \mathbf{s}_t} \frac{1}{t} \sum_{j=1}^t \left(-\mathcal{J}_j(\mathbf{L}\mathbf{s}_j) + \mu_1 \|\mathbf{s}_j\|_1 \right) + \mu_2 \|\mathbf{L}\|_F^2, \quad (1)$$

where $\|\mathbf{v}\|_1$ and $\|\mathbf{Z}\|_F$ denote the L_1 norm of $\mathbf{v} \in \mathbb{R}^p$ and the Frobenius norm of $\mathbf{Z} \in \mathbb{R}^{d \times p}$, respectively. μ_1 and μ_2 denote regularization parameters, and

$\mathcal{J}_j(\mathbf{L}\mathbf{s}_j)$ is the total pay-off for a task j :

$$\mathcal{J}_j(\mathbf{L}\mathbf{s}_j) = \int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}^{(j)}} p_{\mathbf{L}\mathbf{s}_j}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) d\boldsymbol{\tau}^{(j)}.$$

Here, $p_{\mathbf{L}\mathbf{s}_j}^{(j)}(\boldsymbol{\tau}^{(j)})$ and $\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)})$ denote the task-specific trajectory probability and expected total reward:

$$\begin{aligned} p_{\mathbf{L}\mathbf{s}_j}^{(j)}(\boldsymbol{\tau}^{(j)}) &= p_0^{(j)}(\mathbf{x}_0^{(j)}) \prod_{m=0}^{M^{(j)}-1} p^{(j)}(\mathbf{x}_{m+1}^{(j)} | \mathbf{x}_m^{(j)}, \mathbf{u}_m^{(j)}) \pi_{\mathbf{L}\mathbf{s}_j}^{(j)}(\mathbf{u}_m^{(j)} | \mathbf{x}_m^{(j)}) \\ \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) &= \frac{1}{M^{(j)}} \sum_{m=0}^{M^{(j)}} \mathcal{R}^{(j)}(\mathbf{x}_m^{(j)}, \mathbf{u}_m^{(j)}). \end{aligned}$$

As noted in [9], Equation (1) describes batch multi-task learning as opposed to the online lifelong learning setting. This is due to the dependency of the above optimization problem on all trajectories from all tasks. To remedy this problem, the authors employed a local second-order Taylor expansion of the lower-bound to $\mathcal{J}_t(\boldsymbol{\theta}_t)$ around $\tilde{\boldsymbol{\theta}}_t^* = \arg \min_{\tilde{\boldsymbol{\theta}}} -\mathcal{J}_{\mathcal{L},\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}})$, where $\mathcal{J}_{\mathcal{L},\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}) = \text{KL}(p_{\boldsymbol{\theta}}(\boldsymbol{\tau})\mathfrak{R}(\boldsymbol{\tau}) || p_{\tilde{\boldsymbol{\theta}}}(\boldsymbol{\tau}))$, leading to:

$$\min_{\mathbf{L}, \mathbf{s}_1: \mathbf{s}_t} \frac{1}{t} \sum_{j=1}^t \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - \mathbf{L}\mathbf{s}_j \right\|_{\mathbf{H}^{(j)}}^2 + \mu_1 \|\mathbf{s}_j\|_1 \right) + \mu_2 \|\mathbf{L}\|_{\text{F}}^2, \quad (2)$$

where j denotes task j , $\|\mathbf{v}\|_{\mathbf{H}^{(j)}}^2 = \mathbf{v}^\top \mathbf{H}^{(j)} \mathbf{v}$ and $\mathbf{H}^{(j)}$ denotes the Hessian defined as:

$$\mathbf{H}^{(j)} = -\mathbb{E}_{\boldsymbol{\tau}^{(j)} \sim p_{\tilde{\boldsymbol{\theta}}_j^*}(\boldsymbol{\tau}^{(j)})} \left[\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=1}^{M^{(j)}} \nabla_{\tilde{\boldsymbol{\theta}}_j, \tilde{\boldsymbol{\theta}}_j}^2 \log \pi_{\tilde{\boldsymbol{\theta}}_j}(\mathbf{u}_m^{(j)} | \mathbf{s}_m^{(j)}) \right].$$

The problems of PG-ELLA. Although successful, as shown in [9, 23], PG-ELLA is unscalable to a large number of tasks, which is typical for the lifelong learning
105 setting. Here, we differentiate the two inefficiency sources of PG-ELLA. First, we note that although the dependency on all tasks has been remedied using a second-order expansion, computing the operating point $\tilde{\boldsymbol{\theta}}_j^*$ can be expensive. In essence, determining $\tilde{\boldsymbol{\theta}}_j^*$ equates to solving a local MDP (i.e., the j^{th}) defined over $\mathbb{T}^{(j)}$. This led the authors to propose a one-step gradient update given
110 a *small set* of observed trajectories for a task j . Apart from the reduction in

the number of available trajectories, gradient methods are known to exhibit slow convergence rates and are prone to local minima. Consequently, a PG-ELLA update can lead to a low-quality latent repository \mathbf{L} due to the error in approximating $\tilde{\boldsymbol{\theta}}_j^*$.

115 The second inefficiency reducing PG-ELLA’s scalability arises when updating the shared repository \mathbf{L} . On one hand, the update involves an inversion of a $dp \times dp$ matrix leading to a complexity of $\mathcal{O}(d^3p^2)$ in the best case. Consequently, such a method is not scalable for neither high-dimensional policy parametrization nor latent dimensions p . As, as mentioned in [23], gradient-
 120 based methods share the problems mentioned above. On the other hand, as the number of tasks grows large, these centralized methods (updating \mathbf{L} or computing $\tilde{\boldsymbol{\theta}}_j^*$) become intractable due to computational and memory constraints.

5. Scalable Lifelong Policy Search

Our strategy to remedy the two problems above consists of both scaling
 125 PG-ELLA to multiple trajectories and multiple tasks through more processing units, as well as generalizing gradient methods to algorithms with faster convergence rates. We assume the existence of multiple processors defined through undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, with \mathcal{V} denoting the set of nodes and \mathcal{E} the set of edges. We assume a natural ordering among the nodes from $1, \dots, |\mathcal{V}|$ and
 130 $e_{ij} \in \mathcal{E}$ denotes the edge between nodes i and j with $i < j$. We define a block Matrix \mathbf{A} to have $|\mathcal{E}|$ rows and $|\mathcal{V}|$ columns. Suppose edge e_{ij} is the k th edge in \mathcal{E} associating to the k th row in matrix \mathbf{A} , then, the entry (k, i) of matrix \mathbf{A} equals identity matrix $\mathbf{I} \in \mathbb{R}^{d \times d}$, i.e., $\mathbf{A}(k, i) = \mathbf{I}$, and the block $\mathbf{A}(k, j)$ of matrix \mathbf{A} equals $-\mathbf{I}$, i.e., $\mathbf{A}(k, j) = -\mathbf{I}$ and other entries are set to matrix $\mathbf{0}$.
 135 $\mathbf{A}(\mathbf{I})$ denotes the general edge-node incident matrix with identity matrix \mathbf{I} .

The intuition is that we distribute all the tasks among the nodes in the graph and compute the local variables by local variables from other nodes, which connected by edges in the set \mathcal{E} . One potential drawback is the communication cost between nodes.

140 5.1. Scaling the Computation of $\tilde{\theta}_j^*$

As mentioned earlier, the first problem of current lifelong policy search algorithms relates to the scalability in determining the local operating point $\tilde{\theta}_j^*$ for the second-order Taylor expansion in Equation (2). Given a set of trajectories $\mathbb{T}^{(j)}$ for a task j , $\tilde{\theta}_j^*$ is determined as the solution to:

$$\begin{aligned}\tilde{\theta}_j^* &= \arg \min_{\tilde{\theta}} -\mathcal{J}_{\mathcal{L}, \tilde{\theta}}^{(j)}(\tilde{\theta}) = \arg \min_{\tilde{\theta}} \text{KL} \left(p_{\tilde{\theta}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \parallel p_{\tilde{\theta}_j}^{(j)}(\tau^{(j)}) \right) \\ &= \arg \min_{\tilde{\theta}} \int_{\tau^{(j)} \in \mathbb{T}^{(j)}} p_{\tilde{\theta}^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \log \left[\frac{p_{\tilde{\theta}^{(j)}}^{(j)}(\tau^{(j)})}{p_{\tilde{\theta}_j}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)})} \right] d\tau^{(j)}.\end{aligned}$$

Given the above processing unit, we assume a random split of the set $\mathbb{T}^{(j)}$ among the nodes $1, \dots, |\mathcal{V}|$ such that $\mathbb{T}^{(j)} = \cup_{i=1}^{|\mathcal{V}|} \mathbb{T}_i^{(j)}$ and thus rewrite the optimization problem of $\tilde{\theta}_j^*$ as:

$$\min_{\tilde{\theta}_1^{(j)}: \tilde{\theta}_{|\mathcal{V}|}^{(j)}} \sum_{i=1}^{|\mathcal{V}|} \left[\int_{\tau^{(j)} \in \mathbb{T}_i^{(j)}} p_{\tilde{\theta}_i^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \log \left[\frac{p_{\tilde{\theta}_i^{(j)}}^{(j)}(\tau^{(j)})}{p_{\tilde{\theta}_i}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)})} \right] d\tau^{(j)} \right] \quad (3)$$

$$\text{s.t. } \tilde{\theta}_i^{(j)} = \tilde{\theta}_j^{(j)}, \text{ for all } e_{ij} \in \mathcal{E} \text{ and } i < j,$$

where the constraints $\tilde{\theta}_i^{(j)} = \tilde{\theta}_j^{(j)}$, for all $e_{ij} \in \mathcal{E}$ and $i < j$ have been added so as to ensure agreement across the different processing units in \mathcal{G} . To rewrite the constraint above using a compact representation we defined the following: a vector $\theta^{(j)} = [\tilde{\theta}_1^{\text{T},(j)}, \dots, \tilde{\theta}_{|\mathcal{V}|}^{\text{T},(j)}]^{\text{T}} \in \mathbb{R}^{d|\mathcal{V}|}$ and matrix $\tilde{\mathbf{A}} = \mathbf{A}(\mathbf{I}_{d \times d}) \in \mathbb{R}^{d|\mathcal{E}| \times d|\mathcal{V}|}$, where $\mathbf{I}_{d \times d}$ is a d -dimensional identity matrix. Using the above notation, we reformulate the problem in Equation (3) as:

$$\min_{\tilde{\theta}_1^{(j)}: \tilde{\theta}_{|\mathcal{V}|}^{(j)}} \sum_{i=1}^{|\mathcal{V}|} \left[\int_{\tau^{(j)} \in \mathbb{T}_i^{(j)}} p_{\tilde{\theta}_i^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \log \left[\frac{p_{\tilde{\theta}_i^{(j)}}^{(j)}(\tau^{(j)})}{p_{\tilde{\theta}_i}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)})} \right] d\tau^{(j)} \right]$$

$$\text{s.t. } \tilde{\mathbf{A}} \tilde{\theta}^{(j)} = \mathbf{0}.$$

The above is a constraint optimization problem that we aim to solve efficiently. Generally, this problem can be intractable due to the non-convexity of each of the i objectives. Under the broad-class of exponential family policies, however,

each of the local objectives is convex in $\tilde{\boldsymbol{\theta}}_i^{(j)}$ and can thus be solved optimally and efficiently. Our technique for determining the solution relies on a blend between an augmented Lagrangian and a decomposition-coordination procedure similar to the accelerated direction method of multiples ADMM [24] algorithm. In contrast to ADMM, our method is required to determine the feasible point based exclusively on local interactions among the computational nodes. This setting is well known in the distributed optimization literature [25]. Unfortunately, current methods for distributed ADMM mostly focus on the univariate case (as detailed in [25]). The univariate case is limited by its one-dimensional setting — we generalize a distributed version of ADMM to our multi-dimensional setting to propose a scalable and efficient solver for lifelong policy search. Let $\boldsymbol{\lambda} \in \mathbb{R}^{d|\mathcal{E}|}$ be a vector of Lagrange multipliers associating with the constraint $\tilde{\mathbf{A}}\tilde{\boldsymbol{\theta}}^{(j)} = \mathbf{0}$ and a constant $\rho > 0$ of penalty term $\left\| \tilde{\mathbf{A}}\tilde{\boldsymbol{\theta}}^{(j)} \right\|_2^2$. We augmented Lagrangian methods (similar to existing work [24, 25]) so that the Lagrangian function becomes:

$$\begin{aligned} \mathcal{L}_{\text{Aug}}\left(\tilde{\boldsymbol{\theta}}^{(j)}, \boldsymbol{\lambda}\right) &= \sum_{i=1}^{|\mathcal{V}|} \left[\int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}_i^{(j)}} p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}\left(\boldsymbol{\tau}^{(j)}\right) \mathfrak{R}^{(j)}\left(\boldsymbol{\tau}^{(j)}\right) \log \left[\frac{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}\left(\boldsymbol{\tau}^{(j)}\right)}{p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}\left(\boldsymbol{\tau}^{(j)}\right) \mathfrak{R}^{(j)}\left(\boldsymbol{\tau}^{(j)}\right)} \right] d\boldsymbol{\tau}^{(j)} \right] \\ &\quad - \boldsymbol{\lambda}^\top \tilde{\mathbf{A}}\tilde{\boldsymbol{\theta}}^{(j)} + \frac{\rho}{2} \left\| \tilde{\mathbf{A}}\tilde{\boldsymbol{\theta}}^{(j)} \right\|_2^2. \end{aligned}$$

At this stage, we can use standard ADMM for acquiring primal-dual updates and thus determining the solution $\tilde{\boldsymbol{\theta}}_i^{*,(j)}$. Standard ADMM, however, is not readily applicable to our setting due to the need of a distributed solution in which each node performs parameter updates based on only local information exchange between its neighbors. To achieve such a distributed solution, we generalize the approach in [25] to the multi-dimensional setting and associate a dual variable with the constraint on each edge. Each processor i keeps a local decision estimate $\tilde{\boldsymbol{\theta}}_i^{(j)}$ and a vector of dual variables $\boldsymbol{\lambda}_{ki}$ with $k < i$. We define two sets of the neighbors of a node i , denote by \mathbb{P}_i and \mathbb{S}_i . Here, \mathbb{P}_i collects all nodes having an index lower than i , i.e., $\mathbb{P}_i = \{j | e_{ij} \in \mathcal{E}, j < i\}$ and \mathbb{S}_i all

Algorithm 1 Scalable Operating Point Computation

- 1: Initialize $\tilde{\theta}_i^{(j)}$, $\forall i \in \{1, \dots, |\mathcal{V}|\}$, and set $\rho > 0$.
- 2: **for** $k = 0, \dots, K$ **do**
- 3: Each agent i updates, $\tilde{\theta}_i^{(j)}$ in a sequential order from $i = 1, \dots, |\mathcal{V}|$ using:

$$\begin{aligned} \tilde{\theta}_i^{(k+1),(j)} = \arg \min_{\tilde{\theta}_i^{(j)}} & \left[\int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}_i^{(j)}} p_{\tilde{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \log \left[\frac{p_{\tilde{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)})}{p_{\tilde{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)})} \right] d\boldsymbol{\tau}^{(j)} \right. \\ & \left. + \frac{\rho}{2} \sum_{l \in \mathbb{P}_i} \left\| \tilde{\theta}_l^{(k+1),(j)} - \tilde{\theta}_i^{(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right\|_2^2 + \frac{\rho}{2} \sum_{l \in \mathbb{S}_i} \left\| \tilde{\theta}_i^{(j)} - \tilde{\theta}_l^{(k),(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{il}^{(k)} \right\|_2^2 \right]. \end{aligned}$$

- 4: Each agent updates $\boldsymbol{\lambda}_{li}$ for $l \in \mathbb{P}_i$ as:

$$\boldsymbol{\lambda}_{li}^{(k+1)} = \boldsymbol{\lambda}_{li}^{(k)} - \rho \left(\tilde{\theta}_l^{(k+1),(j)} - \tilde{\theta}_i^{(k+1),(j)} \right).$$

5: **end for**

nodes with index higher than i , i.e., $\mathbb{S}_i = \{j | e_{ij} \in \mathcal{E}, i < j\}$. Our algorithm to determine $\tilde{\theta}_j^*$ is summarized by the set of instructions in Algorithm 1. Algorithm 1 consists of two steps. First, primal updates are performed by each node in \mathcal{G} (line 3). Second, given the updated primal, the dual variables are then computed (line 4).

Note that in the case of Gaussian policies, which also can be generalized to any policy from the exponential family, the primal updates on line 3 of Algorithm 1 can be computed in closed form (see Appendix B for derivation):

$$\begin{aligned} \tilde{\theta}_i^{(k+1),(j)} = & \left[-\frac{\rho \text{deg}_i}{2} \mathbf{I}_{d \times d} + \mathbb{E}_{\boldsymbol{\tau}^{(j)} \sim p_{\tilde{\theta}_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=0}^{M^{(j)}-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \right] \right]^{-1} \\ & \times \left[\mathbb{E}_{\boldsymbol{\tau}^{(j)} \sim p_{\tilde{\theta}_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=0}^{M^{(j)}-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \mathbf{u}_m^{(j)} \right] \right. \\ & \left. - \frac{\rho}{2} \left(\sum_{l \in \mathbb{S}_i} \left[\tilde{\theta}_l^{(k),(j)} + \frac{1}{\rho} \boldsymbol{\lambda}_{il}^{(k)} \right] + \sum_{l \in \mathbb{P}_i} \left[\tilde{\theta}_l^{(k+1),(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right] \right) \right]. \end{aligned}$$

where deg_i is the degree of processing unit i , i.e., set of neighbors, and $\Phi(\cdot)$ is a

feature map of the system’s state space. Having the problem solved in computing the linearization operating point, the second step needed for scaling lifelong policy search is handling the inefficiency in updating the shared repository. We
 160 detail the solution to this step in the next section.

Theoretical Guarantees: We now show the theorem derives an improvement to a *linear* convergence rate for our proposed algorithm.

Theorem 1. *The sequence $\{\boldsymbol{\theta}_i^{(k),(j)}, \boldsymbol{\lambda}^{(k)}\}_{k \geq 0}$ consists of*

$$\boldsymbol{\theta}_i^{(k),(j)} = [\tilde{\boldsymbol{\theta}}_1^{\top,(j)}, \dots, \tilde{\boldsymbol{\theta}}_{|\mathcal{V}|}^{\top,(j)}]^\top \text{ and } \boldsymbol{\lambda}^{(k)} = [\boldsymbol{\lambda}_1^{(k),\top}, \dots, \boldsymbol{\lambda}_{|\mathcal{E}|}^{(k),\top}]^\top,$$

where $\boldsymbol{\lambda}_i^{(k)} \in \mathbb{R}^d$ is a dual vector. Following Algorithm 1, the sequence converges linearly with a rate given by $\mathcal{O}(1/k)$.

165 (The interested reader can find the proof in Appendix A.)

5.2. Scaling the Model Update

The lifelong policy search problem is framed as the solution to the optimization problem in Equation (2). It is easy to see that such a problem is not convex in both the shared repository, \mathbf{L} , and the co-efficient vectors, \mathbf{s}_j , simultaneously. Holding one variable fixed, however, leads to a convex problem in the other. Consequently, following such an alternating procedure, one can solve a sequence of convex optimization problems iteratively. The computations over the \mathbf{s}_j vectors are relatively simple to scale. The scalability problem arises when trying to update the latent repository \mathbf{L} due to its dependency on all tasks observed thus far. To scale the computation of the shared repository \mathbf{L} , we follow a similar strategy to that of the previous section. Holding \mathbf{s}_j fixed, computing \mathbf{L} , then we need to solve:

$$\min_{\mathbf{L}} \frac{1}{t} \sum_{j=1}^t \left\| \tilde{\boldsymbol{\theta}}_j^* - \mathbf{L} \mathbf{s}_j \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\mathbf{L}\|_{\mathbb{F}}^2.$$

First, we introduce the $\text{vec}(\mathbf{Z})$ notation, which vectorizes a matrix $\mathbf{Z} \in \mathbb{R}^{d \times p}$ to a vector of size \mathbb{R}^{dp} , i.e.,

$$\text{vec}(\mathbf{Z}) = [a_{1,1}, \dots, a_{d,1}, a_{1,2}, \dots, a_{d,2}, \dots, a_{d,1}, \dots, a_{d,p}]^\top,$$

where $a_{i,j}$ denotes the (i,j) -th element of matrix \mathbf{Z} . Consequently, we can write $\text{vec}(\mathbf{L}\mathbf{s}_j) = (\mathbf{s}_j^\top \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L})$, with \otimes being the Kronecker product and $\mathbf{I}_{d \times d}$ a d -dimensional identity matrix. Hence, we can rewrite the optimization problem in the following equivalent form:

$$\min_{\text{vec}(\mathbf{L})} \frac{1}{t} \sum_{j=1}^t \left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L})\|_2^2,$$

where $\|\mathbf{L}\|_{\mathbb{F}}^2 = \|\text{vec}(\mathbf{L})\|_2^2$. Splitting this equation across \mathcal{G} leads us to the following:

$$\begin{aligned} \min_{\text{vec}(\mathbf{L}_1): \text{vec}(\mathbf{L}_{|\mathcal{V}|})} \sum_{i=1}^{|\mathcal{V}|} \left(\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right) \\ \text{s.t. } \text{vec}(\mathbf{L}_l) = \text{vec}(\mathbf{L}_i) \text{ for all } e_{li} \in \mathcal{E} \text{ and } l < i, \end{aligned}$$

where $\text{vec}(\mathbf{L}_i)$ is the chunk of the shared repository to be learned using processing unit i and t_i is the total number of tasks assigned to that processing unit. Similarly, we can rewrite the constraints compactly in the form of $\tilde{\mathbf{A}}\bar{\text{vec}}(\mathbf{L}) = \mathbf{0}$, where $\bar{\text{vec}}(\mathbf{L}) = [\text{vec}(\mathbf{L}_1)^\top, \dots, \text{vec}(\mathbf{L}_{|\mathcal{V}|})^\top]^\top$. Thus, we can write the augmented Lagrangian as:

$$\begin{aligned} \mathcal{L}_{\text{AugL}}(\bar{\text{vec}}(\mathbf{L}), \boldsymbol{\lambda}_{\mathbf{L}}) = \sum_{i=1}^{|\mathcal{V}|} \left(\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right) \\ - \boldsymbol{\lambda}_{\mathbf{L}}^\top \tilde{\mathbf{A}} \bar{\text{vec}}(\mathbf{L}) + \frac{\rho}{2} \left\| \tilde{\mathbf{A}} \bar{\text{vec}}(\mathbf{L}) \right\|_2^2, \end{aligned}$$

where $\rho > 0$ is the penalty parameter and $\boldsymbol{\lambda}_{\mathbf{L}}$ is a vector of Lagrange multipliers. At this stage, we can follow a similar strategy to that of the previous section to determine the optimal solution $\bar{\text{vec}}(\mathbf{L})^*$. To update the primal variables, each processing unit in \mathcal{G} solves the following problem:

$$\begin{aligned} \text{vec}(\mathbf{L}_i^{(k+1)}) = \arg \min_{\text{vec}(\mathbf{L}_i)} \left[\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right. \\ \left. + \frac{\rho}{2} \sum_{l \in \mathbb{P}_i} \left\| \text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_i) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} \right\|_2^2 + \frac{\rho}{2} \sum_{l \in \mathbb{S}_i} \left\| \text{vec}(\mathbf{L}_i) - \text{vec}(\mathbf{L}_l^{(k)}) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} \right\|_2^2 \right], \end{aligned} \quad (4)$$

Algorithm 2 Scalable Operating Point Computation for \mathbf{L}

- 1: Initialize $\text{vec}(\mathbf{L}_i^{(0)})$, $\forall i \in \{1, \dots, |\mathcal{V}|\}$, and set $\rho > 0$.
- 2: **for** $k = 0, \dots, K$ **do**
- 3: Each agent i updates, $\text{vec}(\mathbf{L}_i^{(k+1)})$ in a sequential order from $i = 1, \dots, |\mathcal{V}|$ using:

$$\begin{aligned} \text{vec}(\mathbf{L}_i^{(k+1)}) = \arg \min_{\text{vec}(\mathbf{L}_i)} & \left[\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right. \\ & + \frac{\rho}{2} \sum_{l \in \mathbb{P}_i} \left\| \text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_i) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} \right\|_2^2 \\ & \left. + \frac{\rho}{2} \sum_{l \in \mathbb{S}_i} \left\| \text{vec}(\mathbf{L}_i) - \text{vec}(\mathbf{L}_l^{(k)}) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},il}^{(k)} \right\|_2^2 \right], \end{aligned}$$

- 4: Each agent updates $\boldsymbol{\lambda}_{li}$ for

$$\boldsymbol{\lambda}_{\mathbf{L},li}^{(k+1)} = \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} - \rho \left(\text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_i^{(k+1)}) \right).$$

- 5: **end for**
-

where k is the iteration number. Given an updated primal, the dual variables, \mathbf{L} , are then computed according to:

$$\boldsymbol{\lambda}_{\mathbf{L},li}^{(k+1)} = \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} - \rho \left(\text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_i^{(k+1)}) \right). \quad (5)$$

Equation (5) is similar to step 4 in Algorithm 1. We provide the pseudocode in Algorithm 2. Note that the update rule can be acquired in a closed form the Gaussian policy assumption, leading to the following: (See Appendix B for details):

$$\begin{aligned} \text{vec}(\mathbf{L}_i^{(k+1)}) = & \left[\frac{1}{t_i} \sum_{j=1}^{t_i} \boldsymbol{\Gamma}_j^\top \mathbf{H}_j \boldsymbol{\Gamma}_j + \frac{\rho \text{deg}_i}{2} \mathbf{I}_{dk \times dk} \right]^{-1} \left(\frac{1}{t_i} \sum_{j=1}^{t_i} \boldsymbol{\Gamma}_j \mathbf{H}_j \tilde{\boldsymbol{\theta}}_j^* \right. \\ & \left. + \frac{\rho}{2} \left(\sum_{l \in \mathbb{S}_i} \left[\text{vec}(\mathbf{L}_l^{(k+1)}) + \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} \right] + \sum_{l \in \mathbb{P}_i} \left[\text{vec}(\mathbf{L}_l^{(k+1)}) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},il}^{(k)} \right] \right) \right), \end{aligned}$$

with $\boldsymbol{\Gamma}_j = \mathbf{s}_j \otimes \mathbf{I}_{d \times d}$. Given the solution to the shared repository, the task-specific coefficients are performed locally using a LASSO [26] step for each of

the tasks t_i .

170 **Theoretical Guarantees:** Our next theorem derives a *linear* convergence rate for our proposed algorithm.

Theorem 2. *The sequence $\left\{ \vec{v}ec(\mathbf{L}^{(k)}), \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} \right\}_{k \geq 0}$ consists of*

$$\vec{v}ec(\mathbf{L}^{(k)}) = \left[\vec{v}ec(\mathbf{L}_1)^{(k), \top}, \dots, \vec{v}ec(\mathbf{L}_{|\mathcal{V}|})^{(k), \top} \right]^{\top} \text{ and } \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} = \left[\boldsymbol{\lambda}_{\mathbf{L}, 1}^{(k), \top}, \dots, \boldsymbol{\lambda}_{\mathbf{L}, |\mathcal{E}|}^{(k), \top} \right]^{\top},$$

where $\boldsymbol{\lambda}_{\mathbf{L}, i}^{(k)} \in \mathbb{R}^{dp}$ is a dual vector. Following Equations (4) and 5, the sequence converges linearly with a rate given by $\mathcal{O}(1/k)$.

(The interested reader can find the proof in Appendix A.)

175 Theorems 1 and 2 show that our technique improves over PG-ELLA in both the policy gradient update step and as well as in the computations of the shared repository. This improvement is significant as it shows linear converges in both steps.

6. Experimental Results

180 Next, we present an empirical validation on five benchmark dynamical systems introduced in [4, 23].

Cart Pole: The cart pole (CP) system is controlled by the cart’s mass m_c in kg , the pole’s mass m_p in kg and the pole’s length l in meters. The state is given by the cart’s position and velocity v , as well as the pole’s angle θ and angular velocity $\dot{\theta}$. The goal is to control the pole in an upright position.

Double Inverted Pendulum: The double inverted pendulum (DIP) is an extension of the cart pole system. It has one cart with mass m_0 in kg and two poles in which the corresponding lengths are l_1 and l_2 in meters. We assume the poles have no mass and there are two masses m_1 and m_2 in kg on the top of each pole. The state consists of the cart’s position x_1 and velocity v_1 , the lower pole’s angle θ_1 and angular velocity $\dot{\theta}_1$, as well as the upper pole’s θ_2 and angular velocity $\dot{\theta}_2$. The goal is also to learn a policy to control the two poles in a specific state.

Helicopter: This linearized model of a CH-47 (HC) tandem-rotor helicopter
 195 assumes horizontal motion at 40 knots. It is characterized by two matrices $A \in \mathbb{R}^{4 \times 4}$ and $B \in \mathbb{R}^{4 \times 2}$. The main goal is to stabilize the helicopter by controlling the collective and differential rotor thrust.

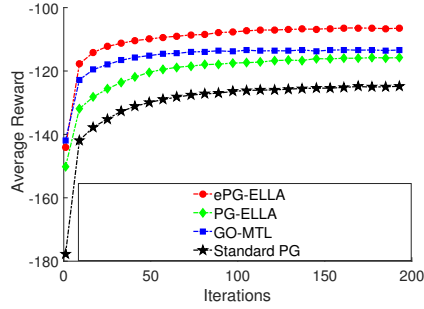
Simple Mass: The simple mass (SM) system is characterized by the spring constant k in N/m , the damping constant d in Ns/m and the mass m in kg .
 200 The system’s state is given by the position x and the velocity, v , of the mass. The goal is to train a policy for guiding the mass to a specific state.

Double Mass: The double mass (DM) is an extension of the simple mass system. It has two masses m_1, m_2 in kg and two springs in which the corresponding springs constants are k_1 and k_2 in N/m , as well as the damping
 205 constant d_1 and d_2 in Ns/m . The state consists of the big mass’ position x_1 and velocity v_1 , as well as the small mass’ position x_2 and velocity v_2 . The goal is also to learn a policy to control the two mass in a specific state.

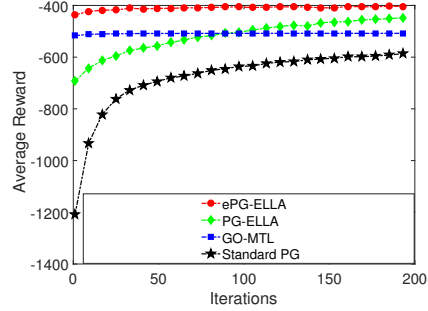
Table 1: Parameter ranges used in the experiments

| CP | DIP | HC |
|---|---|---|
| $m_c, m_p \in [0, 1]$ $l \in [0.2, 0.8]$ | $m_0 \in [1.5, 3.5]$ $m_1, m_2 \in [0.055, 0.1]$ $l_1 \in [0.4, 0.8], l_2 \in [0.5, 0.9]$ | A with rand entry $A(i, j) \in [-32, 3]$ B with rand entry $A(i, j) \in [-9, 1]$ |
| SM | DM | |
| $m \in [3, 5]$ $k \in [1, 7]$ | $m_1 \in [1, 7], m_2 \in [1, 10]$ $k_1 \in [1, 5], k_2 \in [1, 7]$ $d_1, d_2 \in [0.01, 0.1]$ | |

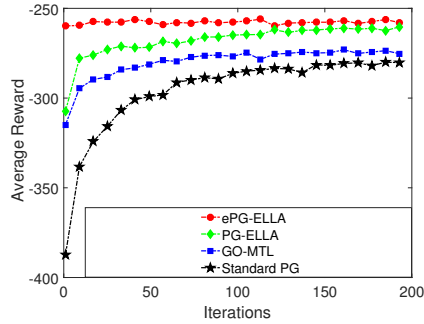
We generated 50 tasks for each domain by varying the dynamical parameters of each of the domains (for 250 in total). For reproducibility, we summarize these
 210 ranges in Table 1. The reward was given by $-\sqrt{\|\mathbf{x}_m - \dot{\mathbf{x}}\|_2^2} - \sqrt{\mathbf{u}_m^\top \mathbf{u}_m}$, where $\dot{\mathbf{x}}$ was the goal state, \mathbf{x}_m is the current state and \mathbf{u}_m is the action that the agent has chosen. We run each task for a total of 200 iterations. At each iteration, the learner observed a task through 50 trajectories of 150 steps and updated



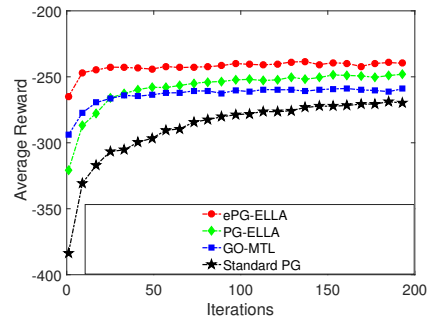
(a) SM



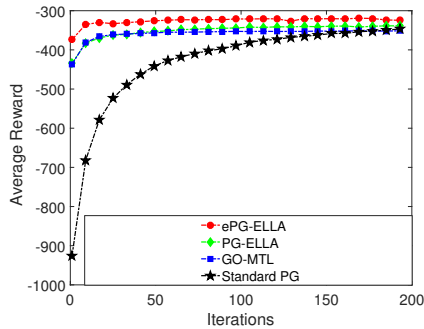
(b) DM



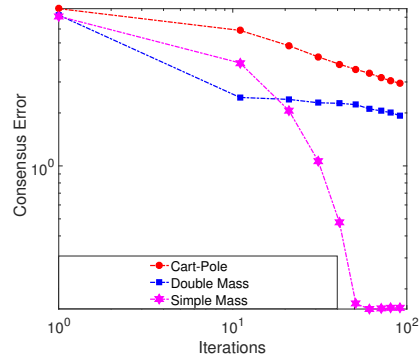
(c) CP



(d) DCP



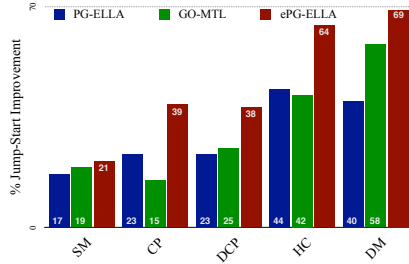
(e) HC



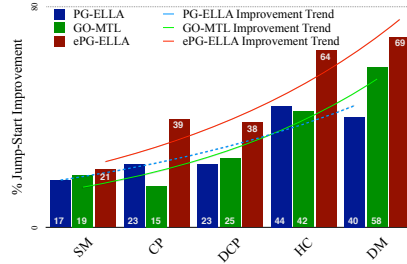
(f) Agreement

Figure 1: Figures (a)-(e) report average reward versus iterations and demonstrate that ePG-ELLA is capable of outperforming other methods. Figure (f) shows the agreement error across processors on 3 sample systems.

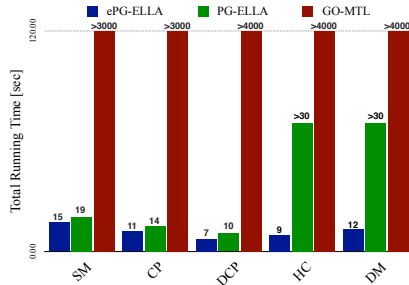
performed algorithmic updates. We used eNAC [27], a standard PG algorithm, as the base learner. We compare our method (ePG-ELLA) to standard PG, PG-



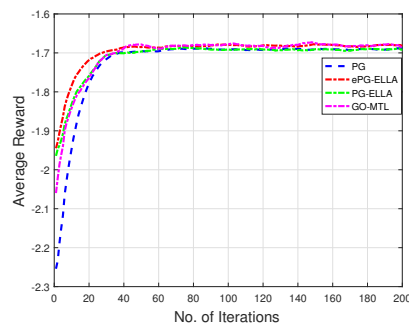
(a) Jump-Start



(b) Asymptotic Results



(c) Time to Convergence



(d) Novel Domain

Figure 2: Figures (a) and (b) show that ePG-ELLA is capable of outperforming others in both the jump-start and asymptotic performance. Figure (c) shows ePG-ELLA takes less time to optimized the objective function. Finally, Figure (d) demonstrates the generalization capability of our method in terms of average reward on unobserved tasks in HC domain.

ELLA [9], and GO-MTL [22]. To distribute our computations, we made use of MATLAB’s parallel pool running on 10 nodes. For distributive computation in ePG-ELLA, we equally assigned 50 tasks to 10 agents, that is, each agent learns 5 tasks in parallel. The edges between agents were generated randomly without overlapping, and the agents are able to exchange information via the graph network. To make sure every node in the graph has at least one predecessor and at least one successor, we connect node 1 to node 2, node 2 to node 3, ..., node 9 and 10. We then randomly assigned 20 edges to all nodes to ensure a connected graph.

Figures 1a—1e report the average reward performance on the different bench-

mark in terms of the number of iterations. In all these systems, our method is capable of outperforming the others in terms of total reward. Figure 1f shows that our method is capable of acquiring feasible solutions by reporting the consensus error on 3 different systems. Figures 2a and 2b demonstrate that our algorithm is capable of outperforming PG-ELLA and GO-MTL in terms of jump-starts and asymptotic performance. Figure 2c demonstrates that our method takes less time to optimize the objective function in terms of time consumption.

Transfer provides significant advantages when the lifelong learning agent faces a novel task domain. To evaluate this, we chose the most complex of the task domains (HC) and trained the lifelong learner by unseen 49 tasks to yield an effectively shared knowledge base for each of the algorithms. Then, we evaluated the algorithms' ability to learn a new (unseen) task from the helicopter domain, comparing the benefits of ePG-ELLA transfer from L_{ePGELLA} with PG-ELLA (from $L_{\text{PG-ELLA}}$), GO-MTL (from $L_{\text{GO-MTL}}$) and PG. Figure 2d depicts the result of learning on a novel domain, averaged over HC tasks, showing that ePG-ELLA converges fastest in this scenario in terms of average reward.

Negative transfer may occur during training with a small probability. However, we run the experiments over 200 iterations and take the average results to smooth the results. Any cases where transfer was negative were statistically overwhelmed by the positive results in our setting.

7. Conclusion and Future Work

We proposed a scalable and efficient lifelong learner for policy gradient methods. Our approach achieves linear convergence rates in both acquiring linearizing operating points and updating shared knowledge bases. We have shown that our new technique can outperform state-of-the-art methods on a variety of benchmark control systems. Future work will include targeting the more-general lifelong learning setting in different domains and running on real world applications such as robotics.

255 **8. Acknowledgements**

We thank the anonymous reviewers for their useful suggestions and comments. This research has taken place in part at the Intelligent Robot Learning (IRL) Lab, Washington State University. IRL's support includes NASA NNX16CD07C, NSF IIS-1149917, NSF IIS-1643614, and USDA 2014-67021-22174.

Appendix A. Convergence and Convergence Rate

In this section, we will show the proof of Theorem 2. Theorem 1 can be proved in a similar way. To show that the sequence $\left\{ \bar{vec}(\mathbf{L}^{(k)}), \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} \right\}_{k \geq 0}$ converges with rate $\mathcal{O}(\frac{1}{k})$, let the function

$$\begin{aligned} \mathcal{F}(\bar{vec}(\mathbf{L})) &= \sum_{i=1}^{|\mathcal{V}|} f_i(\text{vec}(\mathbf{L}_i)) \\ &= \sum_{i=1}^{|\mathcal{V}|} \left(\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right) \end{aligned}$$

and the Lagrangian function

$$\begin{aligned} \mathcal{L}_{\text{Aug}_{\mathbf{L}}}(\bar{vec}(\mathbf{L}), \boldsymbol{\lambda}_{\mathbf{L}}) &= \sum_{i=1}^{|\mathcal{V}|} \left(\frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right) \right) \\ &\quad - \boldsymbol{\lambda}_{\mathbf{L}}^{\top} \tilde{\mathbf{A}} \bar{vec}(\mathbf{L}) + \frac{\rho}{2} \left\| \tilde{\mathbf{A}} \bar{vec}(\mathbf{L}) \right\|_2^2, \end{aligned} \tag{A.1}$$

where

$$\bar{vec}(\mathbf{L}^{(k)}) = \left[\text{vec}(\mathbf{L}_1^{(k)})^{\top}, \dots, \text{vec}(\mathbf{L}_{|\mathcal{V}|}^{(k)})^{\top} \right]^{\top} \quad \text{and} \quad \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} = \left[\boldsymbol{\lambda}_{\mathbf{L},1}^{(k),\top}, \dots, \boldsymbol{\lambda}_{\mathbf{L},|\mathcal{E}|}^{(k),\top} \right]^{\top},$$

with $\boldsymbol{\lambda}_{\mathbf{L},i}^{(k)} \in \mathbb{R}^{dk}$ being a dual vector. It is easy to verify that each cost function

$$f_i(\text{vec}(\mathbf{L}_i)) = \frac{1}{t_i} \sum_{j=1}^{t_i} \left(\left\| \tilde{\boldsymbol{\theta}}_j^* - (\mathbf{s}_j \otimes \mathbf{I}_{d \times d}) \text{vec}(\mathbf{L}_i) \right\|_{\mathbf{H}_j}^2 + \mu_2 \|\text{vec}(\mathbf{L}_i)\|_2^2 \right)$$

is convex. Also, we need the following assumptions

Assumption 1 (Existence of a Saddle Point). *The Lagrangian function $\mathcal{L}_{Aug_{\mathbf{L}}}$ has a saddle point, that is, there is a solution $(\bar{vec}(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*)$ such that*

$$\mathcal{L}_{Aug_{\mathbf{L}}}(\bar{vec}(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}) \leq \mathcal{L}_{Aug_{\mathbf{L}}}(\bar{vec}(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*) \leq \mathcal{L}_{Aug_{\mathbf{L}}}(\bar{vec}(\mathbf{L}), \boldsymbol{\lambda}_{\mathbf{L}}^*) \quad (\text{A.2})$$

for all $\bar{vec}(\mathbf{L}) \in \mathbb{R}^{dp|\mathcal{V}|}$ and $\boldsymbol{\lambda}_{\mathbf{L}} \in \mathbb{R}^{dp}$.

Assumption 2. *The penalty term $\left\| \tilde{\mathbf{A}} \bar{vec}(\mathbf{L}) \right\|_2^2$ is bounded by a constant $\gamma > 0$.*

265 We extended Wei and Ozdaglar's results to multiple dimension [25]. Before we provide the proofs, we need following lemmas:

Lemma 1. *The sequence $\left\{ \bar{vec}(\mathbf{L}^{(k)}), \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} \right\}_{k \geq 0}$ consists of*

$$\bar{vec}(\mathbf{L}^{(k)}) = \left[vec(\mathbf{L}_1^{(k)})^\top, \dots, vec(\mathbf{L}_{|\mathcal{V}|}^{(k)})^\top \right]^\top \quad \text{and} \quad \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} = \left[\boldsymbol{\lambda}_{\mathbf{L},1}^{(k),\top}, \dots, \boldsymbol{\lambda}_{\mathbf{L},|\mathcal{E}|}^{(k),\top} \right]^\top,$$

where $\boldsymbol{\lambda}_{\mathbf{L},i}^{(k)} \in \mathbb{R}^{dp}$ is a dual vector. Then, for all $k \geq 0$, we have following relation,

$$\begin{aligned} & \mathcal{F}(\bar{vec}(\mathbf{L})) - \mathcal{F}(\bar{vec}(\mathbf{L}^{(k+1)})) + \left(\bar{vec}(\mathbf{L}^{(k)}) - \bar{vec}(\mathbf{L}^{(k+1)}) \right)^\top \\ & \cdot \left(-\tilde{\mathbf{A}}^\top \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} - \rho \tilde{\mathbf{A}}^\top \tilde{\mathbf{B}} \left(\bar{vec}(\mathbf{L}^{(k+1)}) - \bar{vec}(\mathbf{L}^{(k)}) \right) + \rho \tilde{\mathbf{B}}^\top \tilde{\mathbf{B}} \left(\bar{vec}(\mathbf{L}^{(k+1)}) - \bar{vec}(\mathbf{L}^{(k)}) \right) \right) \geq 0 \end{aligned}$$

for all $\bar{vec}(\mathbf{L}) \in \mathbb{R}^{dp|\mathcal{V}|}$, where matrix $\tilde{\mathbf{A}}$ is the general edge-node incident matrix and $\tilde{\mathbf{B}} = \min \left\{ \mathbf{0}, \tilde{\mathbf{A}} \right\}$, and the min is taken element-wise.

Proof. Let

$$\begin{aligned} g_i^{(k)}(vec(\mathbf{L}_i)) &= \frac{\rho}{2} \sum_{l \in \mathbb{P}_i} \left\| vec(\mathbf{L}_l^{(k+1)}) - vec(\mathbf{L}_i) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)} \right\|_2^2 \\ &+ \frac{\rho}{2} \sum_{l \in \mathbb{S}_i} \left\| vec(\mathbf{L}_i) - vec(\mathbf{L}_l^{(k)}) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},il}^{(k)} \right\|_2^2 \end{aligned}$$

From equation 4, the optimal value of function $g_i^k + f_i$ is $vec(\mathbf{L}_i^{(k+1)})$, which implies that there exists a subgradient $sg(vec(\mathbf{L}_i^{(k+1)})) \in \partial f_i(vec(\mathbf{L}_i))$ such that $sg(vec(\mathbf{L}_i^{(k+1)})) + \nabla g_i^{(k)}(vec(\mathbf{L}_i^{(k+1)})) = 0$, and we obtain

$$\left[vec(\mathbf{L}_i) - vec(\mathbf{L}_i^{(k+1)}) \right]^\top \left(sg(vec(\mathbf{L}_i^{(k+1)})) + \nabla g_i^{(k)}(vec(\mathbf{L}_i^{(k+1)})) \right) = 0, \quad (\text{A.3})$$

for all $\text{vec}(\mathbf{L}_i) \in \mathbb{R}^{dk}$. Due to the definition of subgradient, we have

$$f_i(\text{vec}(\mathbf{L}_i)) \geq f_i\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) + \left[\text{vec}(\mathbf{L}_i) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right]^\top \cdot \text{sg}\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) \quad (\text{A.4})$$

Combining Equation (A.3) and A.4 yield:

$$f_i(\text{vec}(\mathbf{L}_i)) - f_i\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) + \left[\text{vec}(\mathbf{L}_i) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right]^\top \cdot \nabla g_i^{(k)}\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) \geq 0. \quad (\text{A.5})$$

The gradient of $g_i^{(k)}\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right)$ is

$$\begin{aligned} \nabla g_i^{(k)}\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) &= -\rho \sum_{l \in \mathbb{P}_i} \left(\text{vec}\left(\mathbf{L}_l^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)}\right) \\ &\quad + \rho \sum_{l \in \mathbb{S}_i} \left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_l^{(k)}\right) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},il}^{(k)}\right) \end{aligned}$$

Then, we substitute the gradient in Equation (A.5) with above equation, we have

$$\begin{aligned} &f_i(\text{vec}(\mathbf{L}_i)) - f_i\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) + \left[\text{vec}(\mathbf{L}_i) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right]^\top \\ &\quad \cdot \left(-\rho \sum_{l \in \mathbb{P}_i} \left(\text{vec}\left(\mathbf{L}_l^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k)}\right) \right. \\ &\quad \left. + \rho \sum_{l \in \mathbb{S}_i} \left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_l^{(k)}\right) - \frac{1}{\rho} \boldsymbol{\lambda}_{\mathbf{L},il}^{(k)}\right)\right) \geq 0. \end{aligned}$$

Since the dual variables $\boldsymbol{\lambda}_{\mathbf{L},li}^{(k)}$ have the relation in Equation (5), then,

$$\begin{aligned} &f_i(\text{vec}(\mathbf{L}_i)) - f_i\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) + \left[\text{vec}(\mathbf{L}_i) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right]^\top \\ &\quad \cdot \left(\sum_{l \in \mathbb{P}_i} \boldsymbol{\lambda}_{\mathbf{L},li}^{(k+1)} + \sum_{l \in \mathbb{S}_i} -\boldsymbol{\lambda}_{\mathbf{L},il}^{(k+1)} + \sum_{l \in \mathbb{S}_i} \rho \left(\text{vec}\left(\mathbf{L}_l^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_l^{(k)}\right)\right)\right) \geq 0. \end{aligned}$$

Then, by the definition of matrix $\tilde{\mathbf{A}}$, we can rewrite above equation as

$$\begin{aligned} &f_i(\text{vec}(\mathbf{L}_i)) - f_i\left(\text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right) + \left[\text{vec}(\mathbf{L}_i) - \text{vec}\left(\mathbf{L}_i^{(k+1)}\right)\right]^\top \\ &\quad \cdot \left(-\tilde{\mathbf{A}}(\cdot, i)^\top \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} + \sum_{l \in \mathbb{S}_i} \rho \left(\text{vec}\left(\mathbf{L}_l^{(k+1)}\right) - \text{vec}\left(\mathbf{L}_l^{(k)}\right)\right)\right) \geq 0. \end{aligned}$$

where $\tilde{\mathbf{A}}(\cdot, i)$ denotes the i th column of matrix $\tilde{\mathbf{A}}$. We sum the above equation

from $k = 1$ to $|\mathcal{V}|$,

$$\begin{aligned} & \sum_{i=1}^{|\mathcal{V}|} f_i(\text{vec}(\mathbf{L}_i)) - \sum_{i=1}^{|\mathcal{V}|} f_i(\text{vec}(\mathbf{L}_i^{(k+1)})) + \sum_{i=1}^{|\mathcal{V}|} \left[\text{vec}(\mathbf{L}_i) - \text{vec}(\mathbf{L}_i^{(k+1)}) \right]^\top \\ & \quad \cdot \left(-\tilde{\mathbf{A}}(\cdot, i)^\top \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} + \sum_{l \in \mathbb{S}_i} \rho(\text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_l^{(k)})) \right) \geq 0. \end{aligned} \quad (\text{A.6})$$

By the definition of matrices $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{B}}$, we obtain

$$-\sum_{i=1}^{|\mathcal{V}|} \left[\text{vec}(\mathbf{L}_i) - \text{vec}(\mathbf{L}_i^{(k+1)}) \right]^\top \tilde{\mathbf{A}}(\cdot, i)^\top \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} = -\left[\bar{\text{vec}}(\mathbf{L}) - \bar{\text{vec}}(\mathbf{L}^{(k+1)}) \right]^\top \tilde{\mathbf{A}}^\top \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)}$$

and

$$\begin{aligned} & \sum_{i=1}^{|\mathcal{V}|} \left[\text{vec}(\mathbf{L}_i) - \text{vec}(\mathbf{L}_i^{(k+1)}) \right]^\top \sum_{l \in \mathbb{S}_i} \rho(\text{vec}(\mathbf{L}_l^{(k+1)}) - \text{vec}(\mathbf{L}_l^{(k)})) \\ & = \rho \left[\bar{\text{vec}}(\mathbf{L}) \bar{\text{vec}}(\mathbf{L}^{(k+1)}) \right]^\top \left(-\tilde{\mathbf{A}} + \tilde{\mathbf{B}} \right)^\top \tilde{\mathbf{B}} \left(\bar{\text{vec}}(\mathbf{L}^{(k+1)}) - \bar{\text{vec}}(\mathbf{L}^{(k)}) \right). \end{aligned}$$

The result follows by plugging preceding two equations into Equation (A.6). \square

Lemma 2. *The sequence $\left\{ \bar{\text{vec}}(\mathbf{L}^{(k)}), \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} \right\}_{k \geq 0}$ consists of*

$$\bar{\text{vec}}(\mathbf{L}^{(k)}) = \left[\text{vec}(\mathbf{L}_1^{(k)})^\top, \dots, \text{vec}(\mathbf{L}_{|\mathcal{V}|}^{(k)})^\top \right]^\top \quad \text{and} \quad \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} = \left[\boldsymbol{\lambda}_{\mathbf{L},1}^{(k),\top}, \dots, \boldsymbol{\lambda}_{\mathbf{L},|\mathcal{E}|}^{(k),\top} \right]^\top,$$

where $\boldsymbol{\lambda}_{\mathbf{L},i}^{(k)} \in \mathbb{R}^{d_p}$ is a dual vector. Assume that $(\bar{\text{vec}}(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*)$ is the solution of problem A.1. Let matrix $\tilde{\mathbf{A}}$ is the general edge-node incident matrix and $\tilde{\mathbf{B}} = \min \{ \mathbf{0}, \tilde{\mathbf{A}} \}$, and the min is taken element-wise. Then, we have following relation:

$$\begin{aligned} & \bar{\text{vec}}(\mathbf{L}^{(k+1)})^\top \tilde{\mathbf{A}}^\top (\boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} - \boldsymbol{\lambda}_{\mathbf{L}}^*) + \rho(\bar{\text{vec}}(\mathbf{L}^{(k+1)})^\top \tilde{\mathbf{A}}^\top \tilde{\mathbf{B}} (\bar{\text{vec}}(\mathbf{L}^{(k+1)}) - \bar{\text{vec}}(\mathbf{L}^{(k)}))) \\ & \quad + \rho(\bar{\text{vec}}(\mathbf{L}^*) - \bar{\text{vec}}(\mathbf{L}^{(k+1)}))^\top \tilde{\mathbf{B}}^\top \tilde{\mathbf{B}} (\bar{\text{vec}}(\mathbf{L}^{(k+1)}) - \bar{\text{vec}}(\mathbf{L}^{(k)})) \\ & = \frac{1}{2\rho} \left(\left\| \boldsymbol{\lambda}_{\mathbf{L}}^{(k)} - \boldsymbol{\lambda}_{\mathbf{L}}^* \right\|^2 - \left\| \boldsymbol{\lambda}_{\mathbf{L}}^{(k+1)} - \boldsymbol{\lambda}_{\mathbf{L}}^* \right\|^2 \right) \\ & \quad + \frac{\rho}{2} \left(\left\| \tilde{\mathbf{B}} (\bar{\text{vec}}(\mathbf{L}^{(k)}) - \bar{\text{vec}}(\mathbf{L}^*)) \right\|^2 - \left\| \tilde{\mathbf{B}} (\bar{\text{vec}}(\mathbf{L}^{(k+1)}) - \bar{\text{vec}}(\mathbf{L}^*)) \right\|^2 \right) \\ & \quad - \frac{\rho}{2} \left\| \tilde{\mathbf{B}} (\bar{\text{vec}}(\mathbf{L}^{(k+1)}) - \bar{\text{vec}}(\mathbf{L}^{(k)})) - \tilde{\mathbf{A}} \bar{\text{vec}}(\mathbf{L}^{(k+1)}) \right\|^2 \end{aligned}$$

270 *Proof.* See the proofs of Lemma 4.2 in [25] □

Lemma 3. *Let $(\bar{v}ec(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*)$ be a saddle point of Equation (A.1). Then we have*

$$\tilde{\mathbf{A}}\bar{v}ec(\mathbf{L}^*) = 0.$$

Proof. The result follows by the Assumption 1. □

With aforementioned Lemmas, we are able to show the main theoretical result in the paper.

Proof of Theorem 1. To show the convergence with rate $\mathcal{O}(\frac{1}{k})$, we need an additional variable $y^k = \frac{1}{k} \sum_{s=0}^{k-1} \bar{v}ec(\mathbf{L}^{(s+1)})$ which is the average of $\bar{v}ec(\mathbf{L}^k)$ until time k . Then, We need to bound

$$\mathcal{L}_{\text{Aug}_{\mathbf{L}}}(y^k, \boldsymbol{\lambda}_{\mathbf{L}}^*) - \mathcal{L}_{\text{Aug}_{\mathbf{L}}}(\bar{v}ec(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*)$$

First, $\mathcal{L}_{\text{Aug}_{\mathbf{L}}}(y^k, \boldsymbol{\lambda}_{\mathbf{L}}^*) - \mathcal{L}_{\text{Aug}_{\mathbf{L}}}(\bar{v}ec(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*) \geq 0$ since $(\bar{v}ec(\mathbf{L}^*), \boldsymbol{\lambda}_{\mathbf{L}}^*)$ is the saddle point of Equation (A.1). Then, we derive the upper bound of above equation. By setting $\bar{v}ec(\mathbf{L}) = \bar{v}ec(\mathbf{L}^*)$ in Lemma 1, we have following relation,

$$\begin{aligned} & \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) + (\bar{v}ec(\mathbf{L}^*) - \bar{v}ec(\mathbf{L}^{(s+1)}))^{\top} \\ & \cdot \left(-\tilde{\mathbf{A}}^{\top} \boldsymbol{\lambda}_{\mathbf{L}}^{(s+1)} - \rho \tilde{\mathbf{A}}^{\top} \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) + \rho \tilde{\mathbf{B}}^{\top} \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) \right) \geq 0 \end{aligned}$$

We have $\tilde{\mathbf{A}}\bar{v}ec(\mathbf{L}^*) = 0$ from Lemma 3, so $\bar{v}ec(\mathbf{L}^*)^{\top} \tilde{\mathbf{A}}^{\top} = 0$. Thence, the above equation can be reduced to

$$\begin{aligned} & \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) + \bar{v}ec(\mathbf{L}^{(s+1)})^{\top} \tilde{\mathbf{A}}^{\top} \boldsymbol{\lambda}_{\mathbf{L}}^{(s+1)} \\ & + \left(\rho \bar{v}ec(\mathbf{L}^{(s+1)})^{\top} \tilde{\mathbf{A}}^{\top} \tilde{\mathbf{B}} + \bar{v}ec(\mathbf{L}^{(s)}) + \rho (\bar{v}ec(\mathbf{L}^*) - \bar{v}ec(\mathbf{L}^{(s+1)})) \tilde{\mathbf{B}}^{\top} \tilde{\mathbf{B}} \right) \\ & \cdot (\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) \geq 0 \end{aligned}$$

By adding $(\boldsymbol{\lambda}_{\mathbf{L}}^*)^{\top} \tilde{\mathbf{A}}\bar{v}ec(\mathbf{L}^{(s)}) - (\boldsymbol{\lambda}_{\mathbf{L}}^*)^{\top} \tilde{\mathbf{A}}\bar{v}ec(\mathbf{L}^{(s)})$, we obtain

$$\begin{aligned} & \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) + \bar{v}ec(\mathbf{L}^{(s+1)})^{\top} \tilde{\mathbf{A}}^{\top} \boldsymbol{\lambda}_{\mathbf{L}}^* \\ & + \bar{v}ec(\mathbf{L}^{(s+1)})^{\top} \tilde{\mathbf{A}}^{\top} (\boldsymbol{\lambda}_{\mathbf{L}}^{(s+1)} - \boldsymbol{\lambda}_{\mathbf{L}}^*) + \rho \bar{v}ec(\mathbf{L}^{(s+1)})^{\top} \tilde{\mathbf{A}}^{\top} \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) \\ & + \rho (\bar{v}ec(\mathbf{L}^*) - \bar{v}ec(\mathbf{L}^{(s+1)}))^{\top} \tilde{\mathbf{B}}^{\top} \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) \geq 0 \end{aligned}$$

Using Lemma 2 and Lemma 3, it yields,

$$\begin{aligned}
& \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) + (\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}} \bar{v}ec(\mathbf{L}^{(s+1)}) \\
& + \frac{1}{2\rho} \left(\|\boldsymbol{\lambda}_{\mathbf{L}}^{(s)} - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 - \|\boldsymbol{\lambda}_{\mathbf{L}}^{(s+1)} - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 \right) \\
& + \frac{\rho}{2} \left(\left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^{(s)}) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 - \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^{(s)}) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 \right) \\
& - \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(vec(\mathbf{L}^{(s+1)}) - vec(\mathbf{L}^{(s)})) - \tilde{\mathbf{A}}vec(\mathbf{L}^{(s+1)}) \right\|^2 \geq 0
\end{aligned}$$

Then, we sum aforementioned inequality from $s = 0, \dots, k-1$ and after telescoping cancellation, we have

$$\begin{aligned}
& k\mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - \sum_{s=0}^{k-1} \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) + (\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}}^\top \sum_{s=0}^{k-1} \bar{v}ec(\mathbf{L}^{(s+1)}) \\
& + \frac{1}{2\rho} \|\boldsymbol{\lambda}_{\mathbf{L}}^0 - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^0) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 \\
& \geq \frac{1}{2\rho} \|\boldsymbol{\lambda}_{\mathbf{L}}^k - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^1) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 \\
& + \sum_{s=0}^{k-1} \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^{(s+1)}) - \bar{v}ec(\mathbf{L}^{(s)})) - \tilde{\mathbf{A}}\bar{v}ec(\mathbf{L}^{(s+1)}) \right\|^2 \geq 0
\end{aligned}$$

Since \mathcal{F} is a convex function due to the convexity of f_i , we have $\sum_{s=0}^{k-1} \mathcal{F}(\bar{v}ec(\mathbf{L}^{(s+1)})) \geq k\mathcal{F}(y^k)$ and with the definition of $y^k = \frac{1}{k} \sum_{s=0}^{k-1} \bar{v}ec(\mathbf{L}^{(s+1)})$,

$$\begin{aligned}
& k\mathcal{F}(\bar{v}ec(\mathbf{L}^*)) - k\mathcal{F}(y^k) + k(\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}}^\top y^k \\
& + \frac{1}{2\rho} \|\boldsymbol{\lambda}_{\mathbf{L}}^0 - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^0) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 \geq 0
\end{aligned}$$

Taking $-\frac{1}{k}$ on both sides of the above equation and $(\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}}^\top \bar{v}ec(\mathbf{L}^*) = 0$ by Lemma 3, we obtain

$$\begin{aligned}
& \mathcal{F}(y^k) - (\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}}^\top y^k - \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) + (\boldsymbol{\lambda}_{\mathbf{L}}^*)^\top \tilde{\mathbf{A}}^\top \bar{v}ec(\mathbf{L}^*) \\
& \leq \frac{1}{k} \left(\frac{1}{2\rho} \|\boldsymbol{\lambda}_{\mathbf{L}}^0 - \boldsymbol{\lambda}_{\mathbf{L}}^*\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}}(\bar{v}ec(\mathbf{L}^0) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 \right)
\end{aligned}$$

By Assumption 2 and definition of y^k ,

$$\frac{\rho}{2} \left\| \tilde{\mathbf{A}} y^k \right\|_2^2 \leq \frac{\rho}{2k} \sum_{s=0}^{k-1} \left\| \tilde{\mathbf{A}} \bar{v}ec(\mathbf{L}^{(s+1)}) \right\|_2^2 \leq \frac{\rho}{2k} k\gamma = \frac{\gamma\rho}{2} \quad (\text{A.7})$$

Then,

$$\begin{aligned} & \mathcal{F}(y^k) - (\boldsymbol{\lambda}_L^*)^\top \tilde{\mathbf{A}}^\top y^k + \frac{\rho}{2} \left\| \tilde{\mathbf{A}} y^k \right\|_2^2 y^k - \mathcal{F}(\bar{v}ec(\mathbf{L}^*)) + (\boldsymbol{\lambda}_L^*)^\top \tilde{\mathbf{A}}^\top \bar{v}ec(\mathbf{L}^*) - \frac{\rho}{2} \left\| \tilde{\mathbf{A}} \bar{v}ec(\mathbf{L}^*) \right\|_2^2 \\ & \leq \frac{1}{k} \left(\frac{1}{2\rho} \left\| \boldsymbol{\lambda}_L^0 - \boldsymbol{\lambda}_L^* \right\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^0) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 + \frac{\gamma\rho}{2} \right) \end{aligned} \quad (\text{A.8})$$

Combining Equation (A.8) and (A.7), we have

$$\begin{aligned} & \mathcal{L}_{\text{Aug}_L}(y^k, \boldsymbol{\lambda}_L^*) - \mathcal{L}_{\text{Aug}_L}(\bar{v}ec(\mathbf{L}^*), \boldsymbol{\lambda}_L^*) \\ & \leq \frac{1}{k} \left(\frac{1}{2\rho} \left\| \boldsymbol{\lambda}_L^0 - \boldsymbol{\lambda}_L^* \right\|^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{B}} (\bar{v}ec(\mathbf{L}^0) - \bar{v}ec(\mathbf{L}^*)) \right\|^2 + \frac{\gamma\rho}{2} \right) \end{aligned}$$

Due to the convexity of function \mathcal{F} and linearity of $(\boldsymbol{\lambda}_L^*)^\top \tilde{\mathbf{A}}^\top \bar{v}ec(\mathbf{L}^*)$ and $\frac{\rho}{2} \left\| \tilde{\mathbf{A}} \bar{v}ec(\mathbf{L}^*) \right\|_2^2$, the function $\mathcal{L}_{\text{Aug}_L}(\bar{v}ec(\mathbf{L}), \boldsymbol{\lambda}_L^*)$ is strict convex and has a unique minimizer, i.e., $\bar{v}ec(\mathbf{L}^*)$. Hence if $k \rightarrow \inf$, then the sequence $\left\{ \bar{v}ec(\mathbf{L}^{(k)}), \boldsymbol{\lambda}_L^{(k)} \right\}_{k \geq 0}$ converges to the optimal minimizer $\bar{v}ec(\mathbf{L}^*)$. In practice, we could use the additional variables $y^k = \frac{1}{k} \sum_{s=0}^{k-1} \bar{v}ec(\mathbf{L}^{(s+1)})$ to achieve $\mathcal{O}(\frac{1}{k})$ convergence rate. \square

280 Appendix B. Update Rules

In this section, we provide the derivation of update rules in Section 5.1. The derivation of Section 5.2 can be obtained by similar methods. Using Leibniz integration rule, the gradient can be computed as:

$$\begin{aligned} & \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \left[\int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}_i^{(j)}} p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \log \left[\frac{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)})}{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)})} \right] d\boldsymbol{\tau}^{(j)} \right. \\ & \left. \frac{\rho}{2} \sum_{l \in \mathbb{P}_i} \left\| \tilde{\boldsymbol{\theta}}_l^{(k+1),(j)} - \tilde{\boldsymbol{\theta}}_l^{(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right\|_2^2 + \frac{\rho}{2} \sum_{l \in \mathbb{S}_i} \left\| \tilde{\boldsymbol{\theta}}_l^{(j)} - \tilde{\boldsymbol{\theta}}_l^{(k),(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right\|_2^2 \right] \\ & = \underbrace{\int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}_i^{(j)}} \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \log \left[\frac{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)})}{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)})} \right] d\boldsymbol{\tau}^{(j)}}_{f_{\text{RL}}^{(j)}} \\ & \quad + \rho \text{deg}_i \tilde{\boldsymbol{\theta}}_i^{(j)} - \rho \left(\sum_{l \in \mathbb{S}_i} \left[\tilde{\boldsymbol{\theta}}_l^{(k),(j)} + \frac{1}{\rho} \boldsymbol{\lambda}_{il}^{(k)} \right] + \sum_{l \in \mathbb{P}_i} \left[\tilde{\boldsymbol{\theta}}_l^{(k+1),(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right] \right), \end{aligned}$$

with deg_i is the degree of processing unit i , i.e., set of neighbors. Assuming a Gaussian policy for the i^{th} split of task j , we can $f_{\text{RL}_i^{(j)}}$ as:

$$\begin{aligned} f_{\text{RL}_i^{(j)}} &= p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \log \left[p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \right] - \text{cnst.} \\ &= p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=0}^{M^{(j)}-1} -\frac{1}{2} \left[\mathbf{u}_m^{(j)} - \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right]^\top \boldsymbol{\Sigma}^{-1} \left[\mathbf{u}_m^{(j)} - \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right], \end{aligned}$$

with $\boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \in \mathbb{R}^{m \times d}$ being a feature representation of the state, and $\boldsymbol{\Sigma} \in \mathbb{R}^{m \times m}$ being the covariance matrix of the Gaussian policy. Hence:

$$\begin{aligned} &\int_{\boldsymbol{\tau}^{(j)} \in \mathbb{T}_i^{(j)}} \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \log \left[\frac{p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)})}{p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)})} \right] d\boldsymbol{\tau}^{(j)} \\ &= -\nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \left[\sum_{m=0}^{M^{(j)}-1} \mathbf{u}_m^{(j),\top} \boldsymbol{\Sigma}^{-1} \mathbf{u}_m^{(j)} \right. \\ &\quad \left. - 2 \sum_{m=0}^{M^{(j)}-1} \mathbf{u}_m^{(j),\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right. \\ &\quad \left. + \sum_{m=0}^{M^{(j)}-1} \tilde{\boldsymbol{\theta}}_i^{(j),\top} \boldsymbol{\Phi}^{(j),\top}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right] \\ &= -p_{\boldsymbol{\theta}_i^{(j)}}^{(j)}(\boldsymbol{\tau}^{(j)}) \mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \left[\sum_{m=0}^{M^{(j)}-1} \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \mathbf{u}_m^{(j),\top} \boldsymbol{\Sigma}^{-1} \mathbf{u}_m^{(j)} \right. \\ &\quad \left. - 2 \sum_{m=0}^{M^{(j)}-1} \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \mathbf{u}_m^{(j),\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right. \\ &\quad \left. + \sum_{m=0}^{M^{(j)}-1} \nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \tilde{\boldsymbol{\theta}}_i^{(j),\top} \boldsymbol{\Phi}^{(j),\top}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} \right]. \end{aligned}$$

At this stage, we handle each of the above three gradients separately. The first is zero as the expression does not depend on $\tilde{\boldsymbol{\theta}}_i^{(j)}$. For the second we have:

$$\nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \mathbf{u}_m^{(j),\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} = \boldsymbol{\Phi}^{(j),\top}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \mathbf{u}_m^{(j)}.$$

Furthermore, we have:

$$\nabla_{\tilde{\boldsymbol{\theta}}_i^{(j)}} \tilde{\boldsymbol{\theta}}_i^{(j),\top} \boldsymbol{\Phi}^{(j),\top}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{(j)}(\mathbf{x}_m) \tilde{\boldsymbol{\theta}}_i^{(j)} = (\mathbf{A} + \mathbf{A}^\top) \tilde{\boldsymbol{\theta}}_i^{(j)},$$

with $\mathbf{A} = \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m)$. Since $\mathbf{A} = \mathbf{A}^\top$, then we have:

$$\begin{aligned} \nabla_{\tilde{\theta}_i^{(j)}} \tilde{\theta}_i^{(j),\top} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m) \tilde{\theta}_i^{(j)} &= 2\mathbf{A} \tilde{\theta}_i^{(j)} \\ &= 2\Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m) \tilde{\theta}_i^{(j)}. \end{aligned}$$

Plugging these results back in the original gradient equation, yields:

$$\begin{aligned} &\int_{\tau^{(j)} \in \mathbb{T}_i^{(j)}} \nabla_{\tilde{\theta}_i^{(j)}} p_{\theta_i^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \log \left[\frac{p_{\tilde{\theta}_i^{(j)}}^{(j)}(\tau^{(j)})}{p_{\theta_i^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)})} \right] d\tau^{(j)} \\ &= -p_{\theta_i^{(j)}}^{(j)}(\tau^{(j)}) \mathfrak{R}^{(j)}(\tau^{(j)}) \left[-2 \sum_{m=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \mathbf{u}_m^{(j)} \right. \\ &\quad \left. + 2 \sum_{m=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m) \right] \\ &= 2\mathbb{E}_{\tau^{(j)} \sim p_{\theta_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\tau^{(j)}) \sum_{m=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \mathbf{u}_m^{(j)} \right] \\ &\quad - 2\mathbb{E}_{\tau^{(j)} \sim p_{\theta_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\tau^{(j)}) \sum_{j=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m) \tilde{\theta}_i^{(j)} \right]. \end{aligned}$$

Considering the penalty terms of ADMM, and setting the above to zero, we can update $\tilde{\theta}_i^{(j)}$ in closed-form according to:

$$\begin{aligned} &-2\mathbb{E}_{\tau^{(j)} \sim p_{\theta_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\tau^{(j)}) \sum_{j=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \Phi^{(j)}(\mathbf{x}_m) \tilde{\theta}_i^{(j)} \right] \\ &\quad + 2\mathbb{E}_{\tau^{(j)} \sim p_{\theta_i^{(j)}}^{(j)}(\cdot)} \left[\mathfrak{R}^{(j)}(\tau^{(j)}) \sum_{m=0}^{M^{(j)}-1} \Phi^{(j),\top}(\mathbf{x}_m) \Sigma^{-1} \mathbf{u}_m^{(j)} \right] + \rho \deg_i \tilde{\theta}_i^{(j)} \\ &\quad - \rho \left(\sum_{l \in \mathbb{S}_i} \left[\tilde{\theta}_l^{(k),(j)} + \frac{1}{\rho} \lambda_{il}^{(k)} \right] + \sum_{l \in \mathbb{P}_i} \left[\tilde{\theta}_l^{(k+1),(j)} - \frac{1}{\rho} \lambda_{li}^{(k)} \right] \right) = 0. \end{aligned}$$

Solving the above yields:

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_i^{(k+1),(j)} = & \left[-\frac{\rho \text{deg}_i}{2} \mathbf{I}_{d \times d} + \mathbb{E}_{\boldsymbol{\tau}^{(j)} \sim p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}(\cdot)} \left[\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=0}^{M^{(j)}-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \right] \right]^{-1} \\ & \times \left[\mathbb{E}_{\boldsymbol{\tau}^{(j)} \sim p_{\tilde{\boldsymbol{\theta}}_i^{(j)}}(\cdot)} \left[\mathfrak{R}^{(j)}(\boldsymbol{\tau}^{(j)}) \sum_{m=0}^{M^{(j)}-1} \boldsymbol{\Phi}^{\top,(j)}(\mathbf{x}_m) \boldsymbol{\Sigma}^{-1} \mathbf{u}_m^{(j)} \right] \right. \\ & \left. - \frac{\rho}{2} \left(\sum_{l \in \mathbb{S}_i} \left[\tilde{\boldsymbol{\theta}}_l^{(k),(j)} + \frac{1}{\rho} \boldsymbol{\lambda}_{il}^{(k)} \right] + \sum_{l \in \mathbb{P}_i} \left[\tilde{\boldsymbol{\theta}}_l^{(k+1),(j)} - \frac{1}{\rho} \boldsymbol{\lambda}_{li}^{(k)} \right] \right) \right]. \end{aligned}$$

Similar steps can be performed to determine the update equation for the shared repository by taking the gradient with respect to $\text{vec}(\mathbf{L})$.

- [1] J. Kober, J. R. Peters, Policy search for motor primitives in robotics, in: Advances in neural information processing systems, 2009, pp. 849–856.
- 285 [2] S. A. Murphy, D. W. Oslin, A. J. Rush, J. Zhu, Methodological challenges in constructing effective treatment sequences for chronic psychiatric disorders, Neuropsychopharmacology 32 (2007) 257–262.
- [3] J. Pineau, M. G. Bellemare, A. J. Rush, A. Ghizaru, S. A. Murphy, Constructing evidence-based treatment strategies using methods from computer science, Drug and alcohol dependence 88 (2007) S52–S60.
- 290 [4] R. S. Sutton, A. G. Barto, Introduction to Reinforcement Learning, 1st ed., MIT Press, Cambridge, MA, USA, 1998.
- [5] A. Wilson, A. Fern, S. Ray, P. Tadepalli, Multi-task reinforcement learning: a hierarchical bayesian approach, in: Proceedings of the 24th international conference on Machine learning, ACM, 2007, pp. 1015–1022.
- 295 [6] M. E. Taylor, P. Stone, Transfer Learning for Reinforcement Learning Domains: A Survey, Journal of Machine Learning Research 10 (2009) 1633–1685.
- [7] A. Lazaric, M. Ghavamzadeh, Bayesian multi-task reinforcement learning, in: Proceedings of the 27th International Conference on Machine Learning, 2010.
- 300 [8] H. Li, X. Liao, L. Carin, Multi-task reinforcement learning in partially observable stochastic environments, The Journal of Machine Learning Research 10 (2009) 1131–1186.
- [9] H. Bou-Ammar, E. Eaton, P. Ruvolo, M. Taylor, Online Multi-task Learning for Policy Gradient Methods, in: T. Jebara, E. P. Xing (Eds.), Proceedings of the 31st International Conference on Machine Learning, JMLR Workshop and Conference Proceedings, 2014.
- 305

- [10] R. J. Williams, Simple statistical gradient-following algorithms for connectionist reinforcement learning, *Machine learning* 8 (1992) 229–256.
- [11] S. Bhatnagar, R. S. Sutton, M. Ghavamzadeh, M. Lee, Natural actor–critic algorithms, *Automatica* 45 (2009) 2471–2482.
- [12] J. Peters, S. Schaal, Natural actor-critic, *Neurocomputing* 71 (2008) 1180–1190.
- [13] P. Ruvolo, E. Eaton, Ella: An efficient lifelong learning algorithm, in: *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [14] S. Thrun, J. O’Sullivan, Discovering Structure in Multiple Learning Tasks: The TC Algorithm, in: *International Conference on Machine Learning*, Morgan Kaufmann, 1996.
- [15] W. Caarls, E. Schuitema, Parallel online temporal difference learning for motor control, *IEEE transactions on neural networks and learning systems* 27 (2016) 1457–1468.
- [16] S. Gu, E. Holly, T. Lillicrap, S. Levine, Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates, *arXiv preprint arXiv:1610.00633* (2016).
- [17] A. Yahya, A. Li, M. Kalakrishnan, Y. Chebotar, S. Levine, Collective robot reinforcement learning with distributed asynchronous guided policy search, *arXiv preprint arXiv:1610.00673* (2016).
- [18] S. Levine, C. Finn, T. Darrell, P. Abbeel, End-to-end training of deep visuomotor policies, *Journal of Machine Learning Research* 17 (2016) 1–40.
- [19] M. P. Deisenroth, P. Englert, J. Peters, D. Fox, Multi-task policy search for robotics, in: *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, 2014, pp. 3876–3881.

- [20] A. Wilson, A. Fern, S. Ray, P. Tadepalli, Multi-Task Reinforcement Learning: A Hierarchical Bayesian Approach, in: Machine Learning, Proceedings of the 24th International Conference Corvallis, Oregon, USA, June 20-24, year = 2007,, ????
- 340 [21] M. Snel, S. Whiteson, Learning potential functions and their representations for multi-task reinforcement learning, Autonomous agents and multi-agent systems 28 (2014) 637–681.
- [22] A. Kumar, H. Daume, Learning task grouping and overlap in multi-task learning, in: Proceedings of the 29th International Conference on Machine Learning, 2012, pp. 1383–1390.
- 345 [23] H. Bou Ammar, E. Eaton, J. M. Luna, P. Ruvolo, Autonomous cross-domain knowledge transfer in lifelong policy gradient reinforcement learning, in: Proceedings of the International Joint Conference on Artificial Intelligence, 2015.
- [24] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, Distributed optimization and statistical learning via the alternating direction method of multipliers, Foundations and Trends in Machine Learning 3 (2011) 1–122.
- 350 [25] E. Wei, A. Ozdaglar, Distributed alternating direction method of multipliers, in: Decision and Control, 2012 IEEE 51st Annual Conference on, IEEE, 2012, pp. 5445–5450.
- 355 [26] R. Tibshiranit, Regression shrinkage and selection via the lasso, Journal of the Royal Statistical Society. Series B (Methodological) 58 (1996) pp. 267–288.
- [27] J. Peters, S. Schaal, Natural Actor-Critic, Neurocomputing 71 (2008).