

Thompson sampling for linear quadratic mean-field teams

Mukul Gagrani, Sagar Sudhakara, Aditya Mahajan, Ashutosh Nayyar and Yi Ouyang

Abstract—We consider optimal control of an unknown multi-agent linear quadratic (LQ) system where the dynamics and the cost are coupled across the agents through the mean-field (i.e., empirical mean) of the states and controls. Directly using single-agent LQ learning algorithms in such models results in regret which increases polynomially with the number of agents. We propose a new Thompson sampling based learning algorithm which exploits the structure of the system model and show that the expected Bayesian regret of our proposed algorithm for a system with agents of $|M|$ different types at time horizon T is $\tilde{O}(|M|^{1.5}\sqrt{T})$ irrespective of the total number of agents, where the \tilde{O} notation hides logarithmic factors in T . We present detailed numerical experiments to illustrate the salient features of the proposed algorithm.

I. INTRODUCTION

Linear dynamical systems with a quadratic cost (henceforth referred to as LQ systems) are one of the most commonly used modeling framework in Systems and Control. Part of the appeal of LQ models is that the optimal control action in such models is a linear or affine function of the state; therefore, the optimal policy is easy to identify and easy to implement.

Broadly speaking, the regret of three classes of learning algorithms have been analyzed in the literature: Optimism in the face of uncertainty (OFU) based algorithms, certainty equivalence (CE) based algorithms, and Thompson sampling (TS) based algorithms.

OFU-based algorithms are inspired by the OFU principle for multi-armed bandits [1]. Starting with the work of [2], [3], most of the papers following this approach [4]–[6] provide a high probability bound on regret. As an illustrative example, it is shown in [6] that, with high probability, the regret of a OFU-based learning algorithm is $\tilde{O}(d_x^{0.5}(d_x + d_u)\sqrt{T})$, where d_x is the dimension of the state, d_u is the dimension of the controls, T is the time horizon, and the $\tilde{O}(\cdot)$ notation hides logarithmic terms in T .

Certainty equivalence (CE) is a classical adaptive control algorithm in Systems and Control [7], [8]. Most papers following this approach [9]–[12] also provide a high probability bound on regret. As an illustrative example, it is shown in

[12] that, with high probability, the regret of a CE-based algorithm is $\tilde{O}(d_x^{0.5}d_u\sqrt{T} + d_x^2)$.

Thompson sampling (TS) based algorithms are inspired by TS algorithm for multi-armed bandits [13]. Most papers following this approach [14]–[16] establish a bound on the expected Bayesian regret. As an illustrative example, [15] shows that the regret of a TS-based algorithm is $\tilde{O}(d_x^{0.5}(d_x + d_u)\sqrt{T})$.

Two aspects of these regret bounds are important: the dependence on the time horizon T and the dependence on the dimensions (d_x, d_u) of the state and the controls. For all classes of algorithms mentioned above, the dependence on the time horizon is $\tilde{O}(\sqrt{T})$. Moreover, there are multiple papers which show that, under different assumptions, the regret is lower bounded by $\Omega(\sqrt{T})$ [12], [17]. So, the time dependence in the available regret bounds is nearly order optimal. Similarly, even though the dependence of the regret bound on the dimensions of the state and the control varies slightly for each class of algorithms, [12] recently showed that the regret is lower bounded by $\tilde{\Omega}(d_x^{0.5}d_u\sqrt{T})$. So, there is only a small scope of improvement in the dimension dependence in the regret bounds.

The dependence of the regret bounds on the dimensions of the state and controls is critical for applications such as formation control of robotic swarms and demand response in power grids which have large numbers of agents (which can be of the order of 10^3 to 10^5). In such systems, the effective dimension of the state and the controls is nd_x and nd_u , where n is the number of agents and d_x and d_u are the dimensions of the state and controls of each agent. Therefore, if we take the regret bound of, say, the OFU algorithm proposed in [6], the regret is $\tilde{O}(n^{1.5}d_x^{0.5}(d_x + d_u)\sqrt{T})$. Similar scaling with n holds for CE- and TS-based algorithms. The polynomial dependence on the number of agents is prohibitive and, because of it, the standard regret bounds are of limited value for large-scale systems.

There are many papers in the literature on the design of large-scale systems which exploit some structural property of the system to develop low-complexity design algorithms [18]–[23]. However, there has been very little investigation on the role of such structural properties in developing and analyzing learning algorithms.

Our main contribution is to show that by carefully exploiting the structure of the model, it is possible to design learning algorithms for large-scale LQ systems where the regret does not grow polynomially in the number of agents. In particular, we investigate mean-field coupled control systems, which have emerged as a popular modeling framework in multiple research communities including Control Systems, Economics,

M. Gagrani was with the Department of EE, USC. He is currently with Qualcomm AI Research, San Diego. Email: mgagrani@qti.qualcomm.com
S. Sudhakara and A. Nayyar are with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA. Email: sagarsud@usc.edu; ashutosh@usc.edu.

A. Mahajan, McGill University, Canada. Email: aditya.mahajan@mcgill.ca
Yi Ouyang, Preferred Networks, USA. Email: ouyangyui@gmail.com

The research at USC was supported in part by NSF under grants ECCS 2025732 and ECCS 1750041 and by the Okawa Foundation Research Grant. The work of A. Mahajan was supported in part by the Innovation for Defence Excellence and Security (IDEaS) Program of the Canadian Department of National Defence through grant CFPMN2-30.

Finance, and Statistical Physics [24]–[28]. These models are used in various applications ranging from demand response in smart grids, large scale communication networks, UAVs, financial markets, and many others. We refer the reader to [29] for a survey. There has been considerable interest in reinforcement learning for such models [30]–[35], but all of these papers focus on identifying asymptotically optimal policies and do not characterize regret.

Our main contribution is to design a TS-based algorithm for mean-field teams (which is a specific mean-field model proposed in [22], [23]) and show that the regret scales as $\tilde{O}(|M|^{1.5}d_x^{0.5}(d_x + d_u)\sqrt{T})$, where $|M|$ is the number of types.

We would like to highlight that although we focus on a TS-based algorithm in the paper, it will be clear from the derivation that it is possible to develop OFU- and CE-based algorithms with similar regret bounds. Thus, the main takeaway message of our paper is that there is significant value in developing learning algorithms which exploit the structure of the model.

II. BACKGROUND ON MEAN-FIELD TEAMS

A. Mean-field teams model

We start by describing a slight generalization of the basic model of mean-field teams proposed in [22], [23]. Mean-field teams are also called cooperative mean-field games or mean-field control in the literature [36].

Consider a system with a large population of agents. The agents are heterogeneous and have multiple types. Let $M = \{1, \dots, |M|\}$ denote the set of types of agents, N^m , $m \in M$, denote the set of all agents of type m , and $N = \bigcup_{m \in M} N^m$ denote the set of all agents.

a) States, actions, and their mean-fields: Agents of the same type have the same state and action spaces. In particular, the state and control action of agents of type m take values in $\mathbb{R}^{d_x^m}$ and $\mathbb{R}^{d_u^m}$, respectively. For any generic agent $i \in N^m$ of type m , we use $x_t^i \in \mathbb{R}^{d_x^m}$ and $u_t^i \in \mathbb{R}^{d_u^m}$ to denote its state and control action at time t . We use $\mathbf{x}_t = \text{vec}((x_t^i)_{i \in N})$ and $\mathbf{u}_t = \text{vec}((u_t^i)_{i \in N})$ to denote the global state and control actions of the system at time t .

The empirical mean-field $(\bar{x}_t^m, \bar{u}_t^m)$ of agents of type m , $m \in M$, is defined as the empirical mean of the states and actions of all agents of that type, i.e.,

$$\bar{x}_t^m = \frac{1}{|N^m|} \sum_{i \in N^m} x_t^i \quad \text{and} \quad \bar{u}_t^m = \frac{1}{|N^m|} \sum_{i \in N^m} u_t^i.$$

The empirical mean-field $(\bar{\mathbf{x}}_t, \bar{\mathbf{u}}_t)$ of the entire population is given by

$$\bar{\mathbf{x}}_t = \text{vec}(\bar{x}_t^1, \dots, \bar{x}_t^{|M|}) \quad \text{and} \quad \bar{\mathbf{u}}_t = \text{vec}(\bar{u}_t^1, \dots, \bar{u}_t^{|M|}).$$

As an example, consider the temperature control of a multi-storied office building. In this case, N represents the set of rooms, M represents the set of floors, N^m represents all rooms in floor m , x_t^i represents the temperature in room i , \bar{x}_t^m represents the average temperature in floor m , and $\bar{\mathbf{x}}_t$ represents the collection of average temperature in each floor. Similarly, u_t^i represents the heat exchanged by the air-conditioner

in room i , \bar{u}_t^m represents the average heat exchanged by the air-conditioners in floor m , and $\bar{\mathbf{u}}_t$ represents the collection of average heat exchanged in each floor of the building.

b) System dynamics and per-step cost: The system starts at a random initial state $\mathbf{x}_1 = (x_1^i)_{i \in N}$, whose components are independent across agents. For agent i of type m , the initial state $x_1^i \sim \mathcal{N}(0, \mathbf{X}_1^i)$, and at time $t \geq 1$, the state evolves according to

$$x_{t+1}^i = \mathbf{A}^m x_t^i + \mathbf{B}^m u_t^i + \mathbf{D}^m \bar{\mathbf{x}}_t + \mathbf{E}^m \bar{\mathbf{u}}_t + w_t^i + v_t^m + \mathbf{F}^m v_t^0, \quad (1)$$

where \mathbf{A}^m , \mathbf{B}^m , \mathbf{D}^m , \mathbf{E}^m , \mathbf{F}^m are matrices of appropriate dimensions, $\{w_t^i\}_{t \geq 1}$, $\{v_t^m\}_{t \geq 1}$, and $\{v_t^0\}_{t \geq 1}$ are i.i.d. zero-mean Gaussian processes which are independent of each other and the initial state. In particular, $w_t^i \in \mathbb{R}^{d_x^m}$, $v_t^m \in \mathbb{R}^{d_x^m}$, and $v_t^0 \in \mathbb{R}^{d_v^0}$, and $w_t^i \sim \mathcal{N}(0, \mathbf{W}^i)$, $v_t^m \sim \mathcal{N}(0, \mathbf{V}^m)$, and $v_t^0 \sim \mathcal{N}(0, \mathbf{V}^0)$.

Eq. (1) implies that all agents of type m have similar dynamical couplings. The next state of agent i of type m depends on its current local state and control action, the current mean-field of the states and control actions of the system, and is influenced by three independent noise processes: a local noise process $\{w_t^i\}_{t \geq 1}$, a noise process $\{v_t^m\}_{t \geq 1}$ which is common to all agents of type m , and a global noise process $\{v_t^0\}_{t \geq 1}$ which is common to all agents.

At each time-step, the system incurs a quadratic cost $c(\mathbf{x}_t, \mathbf{u}_t)$ given by

$$c(\mathbf{x}_t, \mathbf{u}_t) = \bar{\mathbf{x}}_t^\top \bar{\mathbf{Q}} \bar{\mathbf{x}}_t + \bar{\mathbf{u}}_t^\top \bar{\mathbf{R}} \bar{\mathbf{u}}_t + \sum_{m \in M} \frac{1}{|N^m|} \sum_{i \in N^m} [(x_t^i)^\top \mathbf{Q}^m x_t^i + (u_t^i)^\top \mathbf{R}^m u_t^i]. \quad (2)$$

Thus, there is a weak coupling in the cost of the agents through the mean-field.

c) Admissible policies and performance criterion: There is a system operator who has access to the states of all agents and control actions and chooses the control action according to a deterministic or randomized policy

$$\mathbf{u}_t = \pi_t(\mathbf{x}_{1:t}, \mathbf{u}_{1:t-1}). \quad (3)$$

Let $\boldsymbol{\theta} = (\boldsymbol{\theta}^m)_{m \in M}$, where $(\boldsymbol{\theta}^m)^\top = [\mathbf{A}^m, \mathbf{B}^m, \mathbf{D}^m, \mathbf{E}^m, \mathbf{F}^m]$, denotes the parameters of the system dynamics. The performance of any policy $\pi = (\pi_1, \pi_2, \dots)$ is given by

$$J(\pi; \boldsymbol{\theta}) = \limsup_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \right]. \quad (4)$$

Let $J(\boldsymbol{\theta})$ to denote the minimum of $J(\pi; \boldsymbol{\theta})$ over all policies.

We are interested in the setup where the system dynamics $\boldsymbol{\theta}$ are unknown and there is a prior p on $\boldsymbol{\theta}$. The Bayesian regret of a policy π operating for a horizon T is defined as

$$R(T; \pi) := \mathbb{E}^\pi \left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) - T J(\boldsymbol{\theta}) \right] \quad (5)$$

where the expectation is with respect to the prior on $\boldsymbol{\theta}$, the noise processes, the initial conditions, and the potential randomizations done by the policy π .

B. Planning solution for mean-field teams

In this section, we summarize the planning solution of mean-field teams presented in [22], [23] for a known system model.

Define the following matrices:

$$\bar{A} = \text{diag}(A^1, \dots, A^{|M|}) + \text{rows}(D^1, \dots, D^{|M|}),$$

$$\bar{B} = \text{diag}(B^1, \dots, B^{|M|}) + \text{rows}(E^1, \dots, E^{|M|}),$$

and let $\bar{Q} = \text{diag}(Q^1, \dots, Q^{|M|}) + \bar{Q}$ and $\bar{R} = \text{diag}(R^1, \dots, R^{|M|}) + \bar{R}$.

It is assumed that the system satisfies the following:

(A1) $\bar{Q} > 0$ and $\bar{R} > 0$. Moreover, for every $m \in M$, $Q^m > 0$ and $R^m > 0$.

(A2) The system (\bar{A}, \bar{B}) is stabilizable.¹ Moreover, for every $m \in M$, the system (A^m, B^m) is stabilizable.

Now, consider the following $|M|+1$ discrete time algebraic Riccati equations (DARE):²

$$\check{S}^m = \text{DARE}(A^m, B^m Q^m, R^m), \quad m \in M, \quad (6a)$$

$$\bar{S} = \text{DARE}(\bar{A}, \bar{B}, \bar{Q}, \bar{R}). \quad (6b)$$

Moreover, define

$$\check{L}^m = -((B^m)^T \check{S}^m B^m + R^m)^{-1} (B^m)^T \check{S}^m A^m, \quad m \in M, \quad (7a)$$

$$\bar{L} = -(\bar{B}^T \bar{S} \bar{B} + \bar{R})^{-1} \bar{B}^T \bar{S} \bar{A}, \quad (7b)$$

and let $\text{rows}(\bar{L}^1, \dots, \bar{L}^{|M|}) = \bar{L}$.

Finally, define $\bar{w}_t^m = \frac{1}{|N^m|} \sum_{i \in N^m} w_t^i$, $\bar{w}_t = \text{vec}(\bar{w}_t^1, \dots, \bar{w}_t^{|M|})$ and $\bar{v}_t = \text{vec}(v_t^1, \dots, v_t^{|M|})$. Let $\check{W}^m = \frac{1}{|N^m|} \sum_{i \in N^m} \text{var}(w_t^i - \bar{w}_t^m)$ and $\check{W} = \text{var}(\bar{w}_t) + \text{diag}(V^1, \dots, V^{|M|}) + \text{diag}(F^1 V^0, \dots, F^{|M|} V^0)$. Note that since the noise processes are i.i.d., these covariances do not depend on time.

Now, split the state x_t^i of agent i of type m into two parts: the *mean-field* state \bar{x}_t^m and the *relative* state $\check{x}_t^i = x_t^i - \bar{x}_t^m$. Do a similar split of the controls: $u_t^i = \bar{u}_t^m + \check{u}_t^i$. Since $\sum_{i \in N^m} \check{x}_t^i = 0$ and $\sum_{i \in N^m} \check{u}_t^i = 0$, the per-step cost (2) can be written as

$$c(x_t, u_t) = \bar{c}(\bar{x}_t, \bar{u}_t) + \sum_{m \in M} \frac{1}{|N^m|} \sum_{i \in N^m} \check{c}^m(\check{x}_t^i, \check{u}_t^i) \quad (8)$$

where $\bar{c}(\bar{x}_t, \bar{u}_t) = \bar{x}_t^T \bar{Q} \bar{x}_t + \bar{u}_t^T \bar{R} \bar{u}_t$ and $\check{c}^m(\check{x}_t^i, \check{u}_t^i) = (\check{x}_t^i)^T Q^m \check{x}_t^i + (\check{u}_t^i)^T R^m \check{u}_t^i$. Moreover, the dynamics of mean-field and the relative components of the state are:

$$\bar{x}_{t+1} = \bar{A} \bar{x}_t + \bar{B} \bar{u}_t + \bar{w}_t + \bar{v}_t + \bar{F} v_t^0 \quad (9)$$

where $\bar{F} = \text{diag}(F^1, \dots, F^{|M|})$ and for any agent i of type m ,

$$\check{x}_t^i = A^m \check{x}_t^i + B^m \check{u}_t^i + \check{w}_t^i, \quad (10)$$

¹System matrices (A, B) are said to be stabilizable if there exists a gain matrix L such that all eigenvalues of $A + BL$ are strictly inside the unit circle.

²For stabilizable (A, B) and $Q > 0$, $\text{DARE}(A, B, Q, R)$ is the unique positive semidefinite solution of $S = A^T S A - (A^T S B)(R + B^T S B)^{-1} (A^T S B) + Q$.

where $\check{w}_t^i = w_t^i - \bar{w}_t^m$.

The result below follows from [23, Theorem 6]³:

Theorem 1 Under assumptions (A1) and (A2), the optimal policy for minimizing the cost (4) is given by

$$u_t^i = \check{L}^m \check{x}_t^i + \bar{L}^m \bar{x}_t. \quad (11)$$

Furthermore, the optimal performance is given by

$$J(\theta) = \sum_{m \in M} \text{Tr}(\check{W}^m \check{S}^m) + \text{Tr}(\bar{W} \bar{S}). \quad (12)$$

a) Interpretation of the planning solution: Note that $\bar{u}_t = \bar{L}_t \bar{x}_t$ is the optimal control for the mean-field system with dynamics (9) and per-step cost $\bar{c}(\bar{x}_t, \bar{u}_t)$. Moreover, for agent i of type m , $\check{u}_t^i = \check{L}_t^m \check{x}_t^i$ is the optimal control for the relative system with dynamics (10) and per-step cost $\check{c}^m(\check{x}_t^i, \check{u}_t^i)$. Theorem 1 shows that at every agent i of type m , we can consider the two decoupled systems—the mean-field system and the relative system—solve them separately, and then simply add their respective controls— \bar{u}_t^m and \check{u}_t^i —to obtain the optimal control action at agent i in the original mean-field team system. We will exploit this feature of the planning solution in order to develop a learning algorithm for mean-field teams.

III. LEARNING FOR MEAN-FIELD TEAMS

For the ease of exposition, we describe the algorithm for the special case when all types are of the same dimension (i.e., $d_x^m = d_x$ and $d_u^m = d_u$ for all $m \in M$) and the same number of agents (i.e., $|N^m| = n$ for all $m \in M$). We further assume that $d_v^0 = d_x$ and $F^m = I$. Moreover, we assume noise covariances are given as $W^i = \sigma_w^2 I$, $i \in N$, $V^m = \sigma_v^2 I$, $m \in M$, and $V^0 = \sigma_{v_0}^2 I$.

The above assumptions are not strictly needed for the analysis but we impose them because, under these assumptions, the covariance matrices $\bar{\Sigma}$ and $\check{\Sigma}^m$ are scaled identity matrices. In particular, for any $m \in M$, $\check{\Sigma}^m = (1 - \frac{1}{n})\sigma_w^2 I =: \check{\sigma}^2 I$ and $\bar{\Sigma} = (\frac{\sigma_w^2}{n} + \sigma_v^2 + \sigma_{v_0}^2)I =: \bar{\sigma}^2 I$. This simpler form of the covariance matrices simplifies the description of the algorithm and the regret bounds.

Following the decomposition presented in Sec. II-B, we define $\bar{\theta}^T = [\bar{A}, \bar{B}]$ to be the parameters of the mean-field dynamics (9) and $(\check{\theta}^m)^T = [A^m, B^m]$ to be the parameters of the relative dynamics (10). We let $\check{S}^m(\check{\theta}^m)$ and $\bar{S}(\bar{\theta})$ denote the solution to the Riccati equations (6) and $\check{L}^m(\check{\theta}^m)$ and $\bar{L}(\bar{\theta})$ denote the corresponding gains (7). Let $\check{J}^m(\check{\theta}^m) = \check{\sigma}^2 \text{Tr}(\check{S}(\check{\theta}^m))$ and $\bar{J}(\bar{\theta}) = \bar{\sigma}^2 \text{Tr}(\bar{S}(\bar{\theta}))$ denote the performance of the m -th relative system and the mean-field system, respectively. As shown in Theorem 1,

$$J(\theta) = \sum_{m \in M} \check{J}^m(\check{\theta}^m) + \bar{J}(\bar{\theta}). \quad (13)$$

³The model considered in [23] did not include common noise, but it is easy to verify that their results continue to hold for models with common noise.

a) *Prior and posterior beliefs*:: We assume that the unknown parameters $\check{\theta}^m$, $m \in M$, lie in compact subsets $\check{\Theta}^m$ of $\mathbb{R}^{(d_x+d_u) \times d_x}$. Similarly, $\bar{\theta}$ lies in a compact subset $\bar{\Theta}$ of $\mathbb{R}^{|M|(d_x+d_u) \times |M|d_x}$. Let $\check{\theta}^m(\ell)$ denote the ℓ -th column of $\check{\theta}^m$. Thus $\check{\theta}^m = \text{cols}(\check{\theta}^m(1), \dots, \check{\theta}^m(d_x))$. Similarly, let $\bar{\theta}(\ell)$ denote the ℓ -th column of $\bar{\theta}$. Thus, $\bar{\theta} = \text{cols}(\bar{\theta}(1), \dots, \bar{\theta}(|M|d_x))$.

We use $\mathcal{N}(\mu, \Sigma)$ to denote the Gaussian distribution with mean μ and covariance Σ and $p|_{\bar{\Theta}}$ to denote the projection of probability distribution p on the set $\bar{\Theta}$.

We assume that the priors \bar{p}_1 and \check{p}_1^m , $m \in M$, on $\bar{\theta}$ and $\check{\theta}^m$, $m \in M$, respectively, satisfy the following:

(A3) \bar{p}_1 is given as:

$$\bar{p}_1(\bar{\theta}) = \left[\prod_{\ell=1}^{|M|d_x} \bar{\lambda}_1^\ell(\bar{\theta}(\ell)) \right] \Big|_{\bar{\Theta}}$$

where for $\ell \in \{1, \dots, |M|d_x\}$, $\bar{\lambda}_1^\ell = \mathcal{N}(\bar{\mu}_1(\ell), \bar{\Sigma}_1)$ with mean $\bar{\mu}_1(\ell) \in \mathbb{R}^{|M|(d_x+d_u)}$ and positive-definite covariance $\bar{\Sigma}_1 \in \mathbb{R}^{|M|(d_x+d_u) \times |M|(d_x+d_u)}$.

(A4) \check{p}_1^m is given as:

$$\check{p}_1^m(\check{\theta}^m) = \left[\prod_{\ell=1}^{d_x} \check{\lambda}_1^{m,\ell}(\check{\theta}^m(\ell)) \right] \Big|_{\check{\Theta}^m}$$

where for $\ell \in \{1, \dots, d_x\}$, $\check{\lambda}_1^{m,\ell} = \mathcal{N}(\check{\mu}_1^m(\ell), \check{\Sigma}_1^m)$ with mean $\check{\mu}_1^m(\ell) \in \mathbb{R}^{d_x+d_u}$ and positive-definite covariance $\check{\Sigma}_1^m \in \mathbb{R}^{(d_x+d_u) \times (d_x+d_u)}$.

These assumptions are similar to the assumptions on the prior in the recent literature on TS for LQ systems [14], [15].

Following the discussion after Theorem 1, we maintain separate posterior distributions on $\bar{\theta}$ and $\check{\theta}^m$, $m \in M$. In particular, we maintain a posterior distribution \bar{p}_t on $\bar{\theta}$ based on the mean-field state and action history as follows: for any Borel subset B of $\mathbb{R}^{|M|(d_x+d_u) \times |M|d_x}$,

$$\bar{p}_t(B) = \mathbb{P}(\bar{\theta} \in B \mid \bar{x}_{1:t}, \bar{u}_{1:t-1}). \quad (14)$$

For every $m \in M$, we also maintain a separate posterior distribution \check{p}_t^m on $\check{\theta}^m$ as follows. At each time $t > 1$, we select an agent $j_{t-1}^m \in N^m$ as $\arg \max_{i \in N^m} (\check{z}_{t-1}^i)^\top \check{\Sigma}_{t-1}^m \check{z}_{t-1}^i$, where $\check{\Sigma}_{t-1}^m$ is a covariance matrix defined recursively by (18b). Then, for any Borel subset B of $\mathbb{R}^{(d_x+d_u) \times d_x}$,

$$\check{p}_t^m(B) = \mathbb{P}(\check{\theta}^m \in B \mid \{\check{x}_s^{j_s^m}, \check{u}_s^{j_s^m}, \check{x}_{s+1}^{j_{s+1}^m}\}_{1 \leq s < t}), \quad (15)$$

See the supplementary file of [37] for a discussion on the rule to select j_{t-1}^m .

For the ease of notation, we use $\bar{z}_t = \text{vec}(\bar{z}_t^1, \dots, \bar{z}_t^{|M|})$, where $\bar{z}_t^m = \text{vec}(\bar{x}_t^m, \bar{u}_t^m)$, and $\check{z}_t^i = \text{vec}(\check{x}_t^i, \check{u}_t^i)$. Then, we can write the dynamics (9)–(10) of the mean-field and the relative systems as

$$\bar{x}_{t+1} = \bar{\theta}^\top \bar{z}_t + \bar{w}_t + \bar{v}_t + v_t^0, \quad (16a)$$

$$\check{x}_{t+1}^i = (\check{\theta}^m)^\top \check{z}_t^i + \check{w}_t^i, \quad \forall i \in N^m, m \in M. \quad (16b)$$

Recall that $\bar{\sigma}^2 = \sigma_w^2/n + \sigma_v^2 + \sigma_{v^0}^2$ and $\check{\sigma}^2 = (1 - \frac{1}{n})\sigma_w^2$.

Lemma 1 *The posterior distributions are as follows:*

1) *The posterior on $\bar{\theta}$ is*

$$\bar{p}_t = \left[\prod_{\ell=1}^{|M|d_x} \bar{\lambda}_t^\ell(\bar{\theta}(\ell)) \right] \Big|_{\bar{\Theta}},$$

where for $\ell \in \{1, \dots, |M|d_x\}$, $\bar{\lambda}_t^\ell = \mathcal{N}(\bar{\mu}_t(\ell), \bar{\Sigma}_t)$, and

$$\bar{\mu}_{t+1}(\ell) = \bar{\mu}_t(\ell) + \frac{\bar{\Sigma}_t \bar{z}_t (\bar{x}_{t+1}(\ell) - \bar{\mu}_t(\ell)^\top \bar{z}_t)}{\bar{\sigma}^2 + (\bar{z}_t)^\top \bar{\Sigma}_t \bar{z}_t}, \quad (17a)$$

$$\bar{\Sigma}_{t+1}^{-1} = \bar{\Sigma}_t^{-1} + \frac{1}{\bar{\sigma}^2} \bar{z}_t \bar{z}_t^\top. \quad (17b)$$

2) *The posterior on $\check{\theta}^m$, $m \in M$, at time t is*

$$\check{p}_t^m(\check{\theta}^m) = \left[\prod_{\ell=1}^{d_x} \check{\lambda}_t^{m,\ell}(\check{\theta}^m(\ell)) \right] \Big|_{\check{\Theta}^m},$$

where for $\ell \in \{1, \dots, d_x\}$, $\check{\lambda}_t^{m,\ell} = \mathcal{N}(\check{\mu}_t^m(\ell), \check{\Sigma}_t^m)$, and

$$\check{\mu}_{t+1}^m(\ell) = \check{\mu}_t^m(\ell) + \frac{\check{\Sigma}_t^m \check{z}_t^{j_t^m} (\check{x}_{t+1}^{j_t^m}(\ell) - \check{\mu}_t^m(\ell)^\top \check{z}_t^{j_t^m})}{\check{\sigma}^2 + (\check{z}_t^{j_t^m})^\top \check{\Sigma}_t^m \check{z}_t^{j_t^m}}, \quad (18a)$$

$$(\check{\Sigma}_{t+1}^m)^{-1} = (\check{\Sigma}_t^m)^{-1} + \frac{1}{\check{\sigma}^2} \check{z}_t^{j_t^m} (\check{z}_t^{j_t^m})^\top. \quad (18b)$$

PROOF Note that the dynamics of \bar{x}_t and \check{x}_t^i in (16) are linear and the noises $\bar{w}_t + \bar{v}_t + v_t^0$ and \check{w}_t^i are Gaussian. Therefore, the result follows from standard results in Gaussian linear regression [38]. ■

b) *The Thompson sampling algorithm*:: We propose a Thompson sampling algorithm referred to as TSDE-MF which is inspired by the TSDE (Thompson sampling with dynamic episodes) algorithm proposed in [14], [15] and the structure of the optimal planning solution for the mean-field teams described in Sec. II-B.

The TSDE-MF algorithm consists of a coordinator \mathcal{C} and $|M| + 1$ actors: a mean-field actor $\bar{\mathcal{A}}$ and a relative actor $\check{\mathcal{A}}^m$, for each $m \in M$. These actors are described below while the whole algorithm is presented in Algorithm 1.

- At each time, the coordinator \mathcal{C} observes the current global state $(x_t^i)_{i \in N}$, computes the mean-field state \bar{x}_t and the relative states $(\check{x}_t^i)_{i \in N}$, and sends the mean-field state \bar{x}_t to be the mean-field actor $\bar{\mathcal{A}}$ and the relative states $\check{x}_t^m = (\check{x}_t^i)_{i \in N^m}$ of the all the agents of type m to the relative actor $\check{\mathcal{A}}^m$. The mean-field actor $\bar{\mathcal{A}}$ computes the mean-field control \bar{u}_t and the relative actor $\check{\mathcal{A}}^m$ computes the relative control $\check{u}_t^m = (\check{u}_t^i)_{i \in N^m}$ (as per the details presented below) and sends it back to the coordinator \mathcal{C} . The coordinator then computes and executes the control action $u_t^i = \bar{u}_t^m + \check{u}_t^i$ for each agent i of type m .
- The mean-field actor $\bar{\mathcal{A}}$ maintains the posterior \bar{p}_t on $\bar{\theta}$ according to (17). The actor works in episodes of dynamic length. Let \bar{t}_k and \bar{T}_k denote the start and the length of episode k , respectively. Episode k ends if the determinant of covariance $\bar{\Sigma}_t$ falls below half of its

value at the beginning of the episode (i.e., $\det(\bar{\Sigma}_t) < 0.5 \det(\bar{\Sigma}_{\bar{t}_k})$) or if the length of the episode is one more than the length of the previous episode (i.e., $t - \bar{t}_k > \bar{T}_{k-1}$). Thus,

$$\bar{t}_{k+1} = \min\{t > \bar{t}_k : \det(\bar{\Sigma}_t) < 0.5 \det(\bar{\Sigma}_{\bar{t}_k}) \text{ or } t - \bar{t}_k > \bar{T}_{k-1}\}. \quad (19)$$

At the beginning of episode k , the mean-field actor $\bar{\mathcal{A}}$ samples a parameter $\bar{\theta}_k$ from the posterior distribution \bar{p}_t . During episode k , the mean-field actor $\bar{\mathcal{A}}$ generates the mean-field controls using the samples $\bar{\theta}_k$, i.e., $\bar{\mathbf{u}}_t = \bar{\mathbf{L}}(\bar{\theta}_k) \bar{\mathbf{x}}_t$.

- Each relative actor $\check{\mathcal{A}}^m$ is similar to the mean-field actor. Actor $\check{\mathcal{A}}^m$ maintains the posterior \check{p}^m on $\check{\theta}^m$ according to (18). The actor works in episodes of dynamic length. The episodes of each relative actor $\check{\mathcal{A}}^m$ and the mean-field actor $\bar{\mathcal{A}}$ are separate from each other.⁴ Let \check{t}_k^m and \check{T}_k^m denote the start and length of episode k , respectively. The termination condition for each episode is similar to that of the mean-field actor $\bar{\mathcal{A}}$. In particular,

$$\check{t}_{k+1}^m = \min\{t > \check{t}_k^m : \det(\check{\Sigma}_t^m) < 0.5 \det(\check{\Sigma}_{\check{t}_k^m}^m) \text{ or } t - \check{t}_k^m > \check{T}_{k-1}^m\}. \quad (20)$$

At the beginning of episode k , the relative actor $\check{\mathcal{A}}^m$ samples a parameter $\check{\theta}_k^m$ from the posterior distribution \check{p}_t^m . During episode k , the relative actor $\check{\mathcal{A}}^m$ generates the relative controls using the sample $\check{\theta}_k^m$, i.e., $\check{\mathbf{u}}_t^m = (\check{\mathbf{L}}^m(\check{\theta}_k^m) \check{\mathbf{x}}_t^m)_{i \in N^m}$.

Note that the algorithm does not depend on the horizon T . A partially distributed version of the algorithm is presented in the conclusion.

c) *Regret bounds*:: We make the following assumption to ensure that the closed loop dynamics of the mean field state and the relative states of each agent are stable. We use the notation $\|\cdot\|$ to denote the induced norm of a matrix.

(A5) There exists $\delta \in (0, 1)$ such that

- for any $\bar{\theta}, \bar{\phi} \in \bar{\Theta}$ where $\bar{\theta}^\top = [\bar{\mathbf{A}}_{\bar{\theta}}, \bar{\mathbf{B}}_{\bar{\theta}}]$, we have $\|\bar{\mathbf{A}}_{\bar{\theta}} + \bar{\mathbf{B}}_{\bar{\theta}} \bar{\mathbf{L}}(\bar{\phi})\| \leq \delta$.
- for any $m \in M$, $\check{\theta}^m, \check{\phi}^m \in \check{\Theta}^m$, where $(\check{\theta}^m)^\top = [\check{\mathbf{A}}_{\check{\theta}^m}, \check{\mathbf{B}}_{\check{\theta}^m}]$, we have $\|\check{\mathbf{A}}_{\check{\theta}^m} + \check{\mathbf{B}}_{\check{\theta}^m} \check{\mathbf{L}}(\check{\phi}^m)\| \leq \delta$.

This assumption is similar to an assumption imposed in the literature on TS for LQ systems [15]. According to Theorem 11 in [12], the assumption is satisfied if

$$\bar{\Theta} = \{(\bar{\mathbf{A}}, \bar{\mathbf{B}}) : \|\bar{\mathbf{A}} - \bar{\mathbf{A}}_0\| \leq \bar{\epsilon}, \|\bar{\mathbf{B}} - \bar{\mathbf{B}}_0\| \leq \bar{\epsilon}\}$$

$$\check{\Theta}^m = \{(\check{\mathbf{A}}^m, \check{\mathbf{B}}^m) : \|\check{\mathbf{A}}^m - \check{\mathbf{A}}_0^m\| \leq \check{\epsilon}^m, \|\check{\mathbf{B}}^m - \check{\mathbf{B}}_0^m\| \leq \check{\epsilon}^m\}$$

for stabilizable $(\bar{\mathbf{A}}_0, \bar{\mathbf{B}}_0)$ and $(\check{\mathbf{A}}_0^m, \check{\mathbf{B}}_0^m)$, and small constants $\bar{\epsilon} \check{\epsilon}^m$ depending on the choice of $(\bar{\mathbf{A}}_0, \bar{\mathbf{B}}_0)$ and $(\check{\mathbf{A}}_0^m, \check{\mathbf{B}}_0^m)$. In other words, the assumption holds when the true system is in a small neighborhood of a known nominal system, and

⁴We use the episode count k as a local variable which is different for each actor.

Algorithm 1 TSDE-MF

```

1: initialize mean-field actor:  $\bar{\Theta}, (\bar{\mu}_1, \bar{\Sigma}_1), \bar{t}_0 = 0, \bar{T}_{-1} = 0, k = 0$ 
2: initialize relative-actor- $m$ :  $\check{\Theta}^m, (\check{\mu}_1^m, \check{\Sigma}_1^m), \check{t}_0^m = 0, \check{T}_{-1}^m = 0, k = 0$ 
3: for  $t = 1, 2, \dots$  do
4:   observe  $(x_t^i)_{i \in N}$ 
5:   compute  $\bar{\mathbf{x}}_t, (\check{\mathbf{x}}_t^m)_{m \in M}$ 
6:    $\bar{\mathbf{u}}_t \leftarrow \text{MEAN-FIELD-ACTOR}(\bar{\mathbf{x}}_t)$ 
7:   for  $m \in M$  do
8:      $\check{\mathbf{u}}_t^m \leftarrow \text{RELATIVE-ACTOR-}m(\check{\mathbf{x}}_t^m)$ 
9:     for  $i \in N^m$  do
10:      agent  $i$  applies control  $u_t^i = \bar{u}_t^m + \check{u}_t^i$ 
11:    end for
12:  end for
13: end for

1: function MEAN-FIELD-ACTOR( $\bar{\mathbf{x}}_t$ )
2:   global var  $t$ 
3:   Update  $\bar{p}_t$  according (17)
4:   if  $t - \bar{t}_k > \bar{T}_{k-1}$  or  $\det(\bar{\Sigma}_t) < 0.5 \det(\bar{\Sigma}_{\bar{t}_k})$  then
5:      $T_k \leftarrow t - \bar{t}_k, k \leftarrow k + 1, \bar{t}_k \leftarrow t$ 
6:     sample  $\bar{\theta}_k \sim \bar{p}_t$ 
7:      $\bar{\mathbf{L}} \leftarrow \bar{\mathbf{L}}(\bar{\theta}_k)$ 
8:   end if
9:   return  $\bar{\mathbf{L}} \bar{\mathbf{x}}_t$ 
10: end function

1: function RELATIVE-ACTOR- $m((\check{\mathbf{x}}_t^i)_{i \in N^m})$ 
2:   global var  $t$ 
3:   Update  $\check{p}_t^m$  according (18)
4:   if  $t - \check{t}_k^m > \check{T}_{k-1}^m$  or  $\det(\check{\Sigma}_t^m) < 0.5 \det(\check{\Sigma}_{\check{t}_k^m}^m)$  then
5:      $T_k^m \leftarrow t - \check{t}_k^m, k \leftarrow k + 1, \check{t}_k^m \leftarrow t$ 
6:     sample  $\check{\theta}_k^m \sim \check{p}_t^m$ 
7:      $\check{\mathbf{L}}^m \leftarrow \check{\mathbf{L}}^m(\check{\theta}_k^m)$ 
8:   end if
9:   return  $(\check{\mathbf{L}}^m \check{\mathbf{x}}_t^i)_{i \in N^m}$ 
10: end function

```

the small neighborhood can be learned with high probability by running some stabilizing procedure [12].

The following result provides an upper bound on the regret of the proposed algorithm.

Theorem 2 Under (A1)–(A5), the regret of TSDE-MF is upper bounded as follows:

$$R(T; \text{TSDE-MF}) \leq \tilde{O}((\bar{\sigma}^2 |M|^{1.5} + \check{\sigma}^2 |M|) d_x^{0.5} (d_x + d_u) \sqrt{T}).$$

Recall that $\bar{\sigma}^2 = \sigma_w^2/n + \sigma_v^2 + \sigma_{v_0}^2$ and $\check{\sigma}^2 = (1 - \frac{1}{n})\sigma_w^2$. So, we can say that $R(T; \text{TSDE-MF}) \leq \tilde{O}(\bar{\sigma}^2 |M|^{1.5} d_x^{0.5} (d_x + d_u) \sqrt{T})$. Compared with the original TSDE regret $\tilde{O}(n^{1.5} |M|^{1.5} \sqrt{T})$ which scales superlinear with the number of agents, the regret of the proposed algorithm is bounded by $\tilde{O}(|M|^{1.5} \sqrt{T})$ irrespective of the total number of agents.

The following special cases are of interest:

- In the absence of common noises (i.e., $\sigma_v^2 = \sigma_{v_0}^2 = 0$), and when $n \gg |M|$, $R(T; \text{TDSE-MF}) \leq \tilde{\mathcal{O}}(\bar{\sigma}^2 |M| d_x^{0.5} (d_x + d_u) \sqrt{T})$.
- For homogeneous systems (i.e., $|M| = 1$), we have $R(T; \text{TDSE-MF}) \leq \tilde{\mathcal{O}}((\bar{\sigma}^2 + \bar{\sigma}^2) d_x^{0.5} (d_x + d_u) \sqrt{T})$. Thus, the scaling with the number of agents is $\tilde{\mathcal{O}}((1 + \frac{1}{n}) \sqrt{T})$.

Note that these results show that in mean-field systems with common noise regret scales as $\mathcal{O}(|M|^{1.5})$ in the number of types, while in mean-field systems without common noise, the regret scales as $\mathcal{O}(|M|)$. Thus, the presence of common noise fundamentally changes the scaling of the learning algorithm.

IV. REGRET ANALYSIS

For the ease of notation, we simply use $R(T)$ instead of $R(T; \text{TSDE-MF})$ in this section. Eq. (13) and (8) imply that the regret may be decomposed as

$$R(T) = \bar{R}(T) + \sum_{m \in M} \frac{1}{n} \sum_{i \in N^m} \check{R}^{i,m}(T) \quad (21)$$

where

$$\begin{aligned} \bar{R}(T) &:= \mathbb{E} \left[\sum_{t=1}^T \bar{c}(\bar{\mathbf{x}}_t, \bar{\mathbf{u}}_t) - T \bar{J}(\bar{\theta}) \right], \\ \check{R}^{i,m}(T) &:= \mathbb{E} \left[\sum_{t=1}^T \check{c}^m(\check{x}_t^i, \check{u}_t^i) - T \check{J}(\check{\theta}^m) \right]. \end{aligned}$$

Note that $\bar{R}(T)$ is the regret associated with the mean-field system and $\check{R}^{i,m}(T)$ is the regret of the i -th relative system of type m . Observe that for the mean-field actor in our algorithm is essentially implementing the TSDE algorithm of [14], [15] for the mean-field system with dynamics (9) and per-step cost $\bar{c}(\bar{\mathbf{x}}_t, \bar{\mathbf{u}}_t)$. This is because:

- 1) As mentioned in the discussion after Theorem 1, we can view $\bar{\mathbf{u}}_t = \bar{\mathbf{L}}(\bar{\theta}) \bar{\mathbf{x}}_t$ as the optimal control action of the mean-field system.
- 2) The posterior distribution \bar{p}_t on $\bar{\theta}$ depends only on $(\bar{\mathbf{x}}_{1:t}, \bar{\mathbf{u}}_{1:t-1})$.

Thus, $\bar{R}(T)$ is precisely the regret of the TSDE algorithm analyzed in [15]. Therefore, we have the following.

Lemma 2 *For the mean-field system,*

$$\bar{R}(T) \leq \tilde{\mathcal{O}}(\bar{\sigma}^2 |M|^{1.5} d_x^{0.5} (d_x + d_u) \sqrt{T}). \quad (22)$$

Unfortunately, we cannot use the same argument to bound $\check{R}^{i,m}(T)$. Even though we can view $\check{u}_t^i = \check{\mathbf{L}}^m(\check{\theta}^m) \check{x}_t^i$ as the optimal control action of the LQ system with dynamics (10), the posterior \check{p}_t^m on $\check{\theta}^m$ depends on terms other than $(\check{x}_{1:t}^i, \check{u}_{1:t-1}^i)$. Therefore, we cannot directly use the results of [15] to bound $\check{R}^{i,m}(T)$. In the rest of this section, we present a bound on $\check{R}^{i,m}(T)$.

For the ease of notation, for any episode k , we use $\check{\mathbf{L}}_k^m$ and $\check{\mathbf{S}}_k^m$ to denote $\check{\mathbf{L}}^m(\check{\theta}_k^m)$ and $\check{\mathbf{S}}^m(\check{\theta}_k^m)$. Recall that the relative value function for average cost LQ problem is $x^\top S x$, where S is the solution to DARE. Therefore, at any time t , episode

k , agent i of type m , and state $\check{x}_t^i \in \mathbb{R}^{d_x}$, with $\check{u}_t^i = \check{\mathbf{L}}_k^m \check{x}_t^i$ and $\check{z}_t^i = \text{vec}(\check{x}_t^i, \check{u}_t^i)$, the average cost Bellman equation is

$$\begin{aligned} \check{J}^m(\check{\theta}_k^m) + (\check{x}_t^i)^\top \check{\mathbf{S}}_k^m \check{x}_t^i &= \check{c}^m(\check{x}_t^i, \check{u}_t^i) \\ &+ \mathbb{E}[(\check{\theta}_k^m)^\top \check{z}_t^i + \check{w}_t^i]^\top \check{\mathbf{S}}_k^m ((\check{\theta}_k^m)^\top \check{z}_t^i + \check{w}_t^i). \end{aligned}$$

Adding and subtracting $\mathbb{E}[(\check{x}_{t+1}^i)^\top \check{\mathbf{S}}_k^m \check{x}_{t+1}^i \mid \check{z}_t^i]$ and noting that $\check{x}_{t+1}^i = (\check{\theta}^m)^\top \check{z}_t^i + \check{w}_t^i$, we get that

$$\begin{aligned} \check{c}^m(\check{x}_t^i, \check{u}_t^i) &= \check{J}^m(\check{\theta}_k^m) + (\check{x}_t^i)^\top \check{\mathbf{S}}_k^m \check{x}_t^i - \mathbb{E}[(\check{x}_{t+1}^i)^\top \check{\mathbf{S}}_k^m \check{x}_{t+1}^i \mid \check{z}_t^i] \\ &+ ((\check{\theta}^m)^\top \check{z}_t^i)^\top \check{\mathbf{S}}_k^m ((\check{\theta}^m)^\top \check{z}_t^i) - ((\check{\theta}_k^m)^\top \check{z}_t^i)^\top \check{\mathbf{S}}_k^m ((\check{\theta}_k^m)^\top \check{z}_t^i). \end{aligned} \quad (23)$$

Let \check{K}_T^m denote the number of episodes of the relative systems of type m until the horizon T . For each $k > \check{K}_T^m$, we define \check{t}_k^m to be $T+1$. Then, using (23), we have that for any agent i of type m ,

$$\begin{aligned} \check{R}^{i,m}(T) &= \mathbb{E} \left[\underbrace{\sum_{k=1}^{\check{K}_T^m} \check{J}_k^m \check{J}^m(\check{\theta}_k^m) - T \check{J}^m(\check{\theta}^m)}_{\text{regret due to sampling error} =: \check{R}_0^{i,m}(T)} \right] \\ &+ \mathbb{E} \left[\underbrace{\sum_{k=1}^{\check{K}_T^m} \sum_{t=\check{t}_k^m}^{\check{t}_{k+1}^m-1} [(\check{x}_t^i)^\top \check{\mathbf{S}}_k^m \check{x}_t^i - (\check{x}_{t+1}^i)^\top \check{\mathbf{S}}_k^m \check{x}_{t+1}^i]}_{\text{regret due to time-varying controller} =: \check{R}_1^{i,m}(T)} \right] \\ &+ \mathbb{E} \left[\underbrace{\sum_{k=1}^{\check{K}_T^m} \sum_{t=\check{t}_k^m}^{\check{t}_{k+1}^m-1} [((\check{\theta}^m)^\top \check{z}_t^i)^\top \check{\mathbf{S}}_k^m ((\check{\theta}^m)^\top \check{z}_t^i) - ((\check{\theta}_k^m)^\top \check{z}_t^i)^\top \check{\mathbf{S}}_k^m ((\check{\theta}_k^m)^\top \check{z}_t^i)]}_{\text{regret due to model mismatch} =: \check{R}_2^{i,m}(T)} \right]. \end{aligned} \quad (24)$$

Lemma 3 *The terms in (24) are bounded as follows:*

- 1) $\check{R}_0^{i,m}(T) \leq \tilde{\mathcal{O}}(\bar{\sigma}^2 \sqrt{(d_x + d_u)T})$.
- 2) $\check{R}_1^{i,m}(T) \leq \tilde{\mathcal{O}}(\bar{\sigma}^2 \sqrt{(d_x + d_u)T})$.
- 3) $\check{R}_2^{i,m}(T) \leq \tilde{\mathcal{O}}(\bar{\sigma}^2 (d_x + d_u) \sqrt{d_x T})$.

PROOF We provide an outline of the proof. See the supplementary file of [37] for complete details.

The first term $\check{R}_0^{i,m}(T)$ can be bounded using the basic property of Thompson sampling: for any measurable function f , $\mathbb{E}[f(\check{\theta}_k^m)] = \mathbb{E}[f(\check{\theta}^m)]$ because $\check{\theta}_k^m$ is a sample from the posterior distribution on $\check{\theta}^m$.

Note that the second term $\check{R}_1^{i,m}(T)$ is a telescopic sum, which we can simplify to establish

$$\check{R}_1^{i,m}(T) \leq \mathcal{O}(\mathbb{E}[\check{K}_T^m (\check{X}_T^i)^2]),$$

where $\check{X}_t^i = \max_{1 \leq t \leq T} \|\check{x}_t^i\|$ is the maximum norm of the relative state along the entire trajectory. The final bound on $\check{R}_1^{i,m}(T)$ can be obtained by bounding \check{K}_T^m and $\mathbb{E}[(\check{X}_T^i)^2]$.

Using the sampling condition for \check{p}_t^m and an existing bound in the literature, we first establish that

$$\check{R}_2^{i,m}(T) \leq \sqrt{\mathbb{E}[(\check{X}_T^i)^2 \sum_{t=1}^T (\check{z}_t^i)^\top \check{\Sigma}_t^m \check{z}_t^i]} \times \tilde{\mathcal{O}}(\sqrt{T})$$

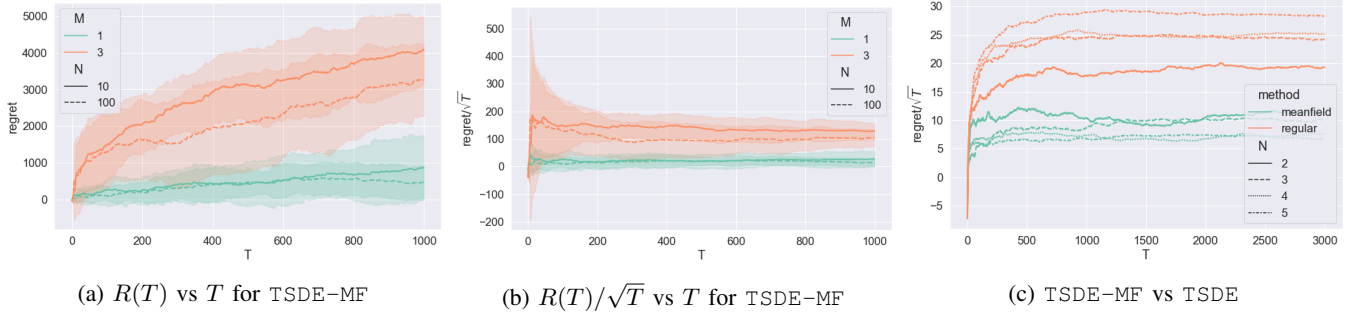


Fig. 1: Expected regret vs time.

Then, we upper bound $(\check{z}_t^i)^\top \check{\Sigma}_t^m \check{z}_t^i$ by $(\check{z}_t^m)^\top \check{\Sigma}_t^m \check{z}_t^m$ which follows from the definition of \check{j}_t^m . Finally, we show that $\mathbb{E}[(\check{X}_T^i)^2 \sum_t (\check{z}_t^i)^\top \check{\Sigma}_t^m \check{z}_t^m]$ is $\tilde{O}(1)$ using the fact that $(\check{\Sigma}_t^m)^{-1}$ is obtained by linearly combining $\{\check{z}_s^m (\check{z}_s^m)^\top\}_{1 \leq s \leq t}$ as in (18b). ■

Combining the three bounds in Lemma 3, we get that

$$\check{R}^{i,m}(T) \leq \tilde{O}(\check{\sigma}^2 d_x^{0.5} (d_x + d_u) \sqrt{T}). \quad (25)$$

By substituting (22) and (25) in (21), we get the result of Theorem 2.

V. NUMERICAL EXPERIMENTS

In this section, we illustrate the performance of TSDE-MF for a homogeneous (i.e., $|M| = 1$) mean-field LQ system for different values of the number n of agents, with $A = 1, B = 0.3, D = 0.5, E = 0.2, Q = 1, \bar{Q} = 1, R = 1$, and $\bar{R} = 0.5$. We set the local noise variance $\sigma_w^2 = 1$.

For the regret plots in Figure 1a,1b, we set the common noise variance to $\sigma_v^2 + \sigma_{v_0}^2 = 1$. The prior distribution used in the simulation are set according to (A3) and (A4) with $\check{\mu}(\ell) = [1, 1]$, $\bar{\mu}(\ell) = [1, 1]$, $\bar{\Sigma}_1 = I$, and $\check{\Sigma}_1 = I$, $\check{\Theta} = \{\check{\theta} : A + B\bar{L}(\check{\theta}) \leq \delta\}$, $\bar{\Theta} = \{\bar{\theta} : A + D + (B + E)\bar{L}(\bar{\theta}) \leq \delta\}$ and $\delta = 0.99$.

In the comparison of TSDE-MF method with TSDE in Figure 1c, we consider the same dynamics and cost parameters as above but without common noise (i.e. $\sigma_v^2 + \sigma_{v_0}^2 = 0$).

a) *Empirical evaluation of regret*:: We run the system for 500 different sample paths and plot the mean and standard deviation of the expected regret $R(T)$ for $T = 5000$. The regret for different values of n is shown in 1a–1b. As seen from the plots, the regret reduces with the number of agents and $R(T)/\sqrt{T}$ converges to a constant. Thus, the empirical regret matches the upper bound of $\tilde{O}((1 + \frac{1}{n})\sqrt{T})$ obtained in Theorem 2.

b) *Comparison with naive TSDE algorithm*:: We compare the performance of TSDE-MF with that of directly using the TSDE algorithm presented in [14], [15] for different values of n . The results are shown in Fig. 1c. As seen from the plots, the regret of TSDE-MF is smaller than TSDE but more importantly, the regret of TSDE-MF reduces with n while that of TSDE increases with n . This matches their respective upper bounds of $\tilde{O}((1 + \frac{1}{n})\sqrt{T})$ and $\tilde{O}(n^{1.5}\sqrt{T})$. These plots clearly illustrate the significance of our results even for small values of n .

VI. CONCLUSION

We consider the problem of controlling an unknown LQ mean-field team. The planning solution (i.e., when the model is known) for mean-field teams is obtained by solving the mean-field system and the relative systems separately. Inspired by this feature, we propose a TS-based learning algorithm TSDE-MF which separately tracks the parameters θ and $\bar{\theta}^m$ of the mean-field and the relative systems, respectively. The part of the TSDE-MF algorithm that learns the mean-field system is similar to the TSDE algorithm for single agent LQ systems proposed in [14], [15] and its regret can be bounded using the results of [14], [15]. However, the part of the TSDE-MF algorithm that learns the relative component is different and we cannot directly use the results of [14], [15] to bound its regret. Our main technical contribution is to provide a bound on the regret on the relative system, which allows us to bound the total regret under TSDE-MF.

a) *Distributed implementation of the algorithm*:: It is possible to implement Algorithm 1 in a distributed manner as follows. Instead of a centralized coordinator which collects all the observations and computes all the controls, we can consider an alternative implementation in which there is an actor \mathcal{A}^m associated with type m and a mean-field actor $\bar{\mathcal{A}}$. Each agent observes its local state and action. The actor \mathcal{A}^m for type m computes (j_t^m, \bar{x}_t^m) using a distributed algorithm, sends \bar{x}_t^m to the mean-field actor, and locally computes $\bar{L}^m(\bar{\theta}_k)$. The mean-field actor computes $\bar{L}(\bar{\theta}_k)$ and sends the m -th block column $\bar{L}^m(\bar{\theta}_k)$ to actors \mathcal{A}^m . Each actor \mathcal{A}^m then sends $(\bar{x}_t^m, \bar{L}^m(\bar{\theta}_k), \bar{L}^m(\bar{\theta}_k))$ to each agent of type m using a distributed algorithm. Each agent then applies the control law (11).

REFERENCES

- [1] P. Auer, N. Cesa-Bianchi, and P. Fischer, “Finite-time analysis of the multiarmed bandit problem,” *Machine learning*, vol. 47, no. 2-3, pp. 235–256, 2002.
- [2] M. C. Campi and P. Kumar, “Adaptive linear quadratic Gaussian control: the cost-biased approach revisited,” *SIAM Journal on Control and Optimization*, vol. 36, no. 6, pp. 1890–1907, 1998.
- [3] Y. Abbasi-Yadkori and C. Szepesvári, “Regret bounds for the adaptive control of linear quadratic systems,” in *Annual Conference on Learning Theory*, pp. 1–26, 2011.
- [4] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, “Finite time analysis of optimal adaptive policies for linear-quadratic systems,” arXiv:1711.07230, 2017.

- [5] A. Cohen, T. Koren, and Y. Mansour, "Learning linear-quadratic regulators efficiently with only \sqrt{T} regret," in *International Conference on Machine Learning*, pp. 1300–1309, PMLR, 2019.
- [6] M. Abeille and A. Lazaric, "Efficient optimistic exploration in linear-quadratic regulators via lagrangian relaxation," in *International Conference on Machine Learning*, pp. 23–31, PMLR, 2020.
- [7] G. Goodwin, P. Ramadge, and P. Caines, "Discrete time stochastic multivariable adaptive control," *IEEE Transactions on Automatic Control*, vol. 19, pp. 449–456, June 1980.
- [8] G. Goodwin, P. Ramadge, and P. Caines, "Discrete time stochastic adaptive control," *SIAM J. Control and Optimization*, vol. 19, pp. 829–853, Nov. 1981.
- [9] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "Regret bounds for robust adaptive control of the linear quadratic regulator," in *Neural Information Processing Systems*, pp. 4192–4201, 2018.
- [10] H. Mania, S. Tu, and B. Recht, "Certainty equivalent control of LQR is efficient," arXiv:1902.07826, 2019.
- [11] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, "Input perturbations for adaptive control and learning," *Automatica*, vol. 117, p. 108950, 2020.
- [12] M. Simchowitz and D. Foster, "Naive exploration is optimal for online lqr," in *International Conference on Machine Learning*, pp. 8937–8948, PMLR, 2020.
- [13] S. Agrawal and N. Goyal, "Analysis of thompson sampling for the multi-armed bandit problem," in *Conference on Learning Theory*, 2012.
- [14] Y. Ouyang, M. Gagrani, and R. Jain, "Control of unknown linear systems with thompson sampling," in *Allerton Conference on Communication, Control, and Computing*, pp. 1198–1205, 2017.
- [15] Y. Ouyang, M. Gagrani, and R. Jain, "Posterior sampling-based reinforcement learning for control of unknown linear systems," *IEEE Transactions on Automatic Control*, 2019.
- [16] M. Abeille and A. Lazaric, "Improved regret bounds for thompson sampling in linear quadratic control problems," in *International Conference on Machine Learning*, pp. 1–9, 2018.
- [17] A. Cassel, A. Cohen, and T. Koren, "Logarithmic regret for learning linear quadratic regulators efficiently," in *International Conference on Machine Learning*, pp. 1328–1337, PMLR, 2020.
- [18] J. Lunze, "Dynamics of strongly coupled symmetric composite systems," *International Journal of Control*, vol. 44, no. 6, pp. 1617–1640, 1986.
- [19] M. K. Sundareshan and R. M. Elbanna, "Qualitative analysis and decentralized controller synthesis for a class of large-scale systems with symmetrically interconnected subsystems," *Automatica*, vol. 27, no. 2, pp. 383–388, 1991.
- [20] G.-H. Yang and S.-Y. Zhang, "Structural properties of large-scale systems possessing similar structures," *Automatica*, vol. 31, no. 7, pp. 1011–1017, 1995.
- [21] S. C. Hamilton and M. E. Broucke, "Patterned linear systems," *Automatica*, vol. 48, no. 2, pp. 263–272, 2012.
- [22] J. Arabneydi and A. Mahajan, "Team-optimal solution of finite number of mean-field coupled lqg subsystems," in *Conf. Decision and Control*, (Kyoto, Japan), Dec. 2015.
- [23] J. Arabneydi and A. Mahajan, "Linear Quadratic Mean Field Teams: Optimal and Approximately Optimal Decentralized Solutions," 2016. arXiv:1609.00056.
- [24] J.-M. Lasry and P.-L. Lions, "Mean field games," *Japanese Journal of Mathematics*, vol. 2, no. 1, pp. 229–260, 2007.
- [25] M. Huang, P. E. Caines, and R. P. Malhamé, "Large-population cost-coupled LQG problems with nonuniform agents: individual-mass behavior and decentralized epsilon-Nash equilibria," *IEEE Transactions on Automatic Control*, vol. 52, no. 9, pp. 1560–1571, 2007.
- [26] M. Huang, P. E. Caines, and R. P. Malhamé, "Social optima in mean field LQG control: centralized and decentralized strategies," *IEEE Transactions on Automatic Control*, vol. 57, no. 7, pp. 1736–1751, 2012.
- [27] G. Y. Weintraub, C. L. Benkard, and B. V. Roy, "Oblivious Equilibrium: A Mean Field Approximation for Large-Scale Dynamic Games," in *Neural Information Processing Systems*, pp. 1489–1496, Dec. 2005.
- [28] G. Y. Weintraub, C. L. Benkard, and B. Van Roy, "Markov perfect industry dynamics with many firms," *Econometrica*, vol. 76, no. 6, pp. 1375–1411, 2008.
- [29] D. A. Gomes and J. Saúde, "Mean field games models—a brief survey," *Dynamic Games and Applications*, vol. 4, no. 2, pp. 110–154, 2014.
- [30] Y. Yang, R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang, "Mean field multi-agent reinforcement learning," in *International Conference on Machine Learning*, pp. 5567–5576, Jul 2018.
- [31] J. Subramanian and A. Mahajan, "Reinforcement learning in stationary mean-field games," in *International Conference on Autonomous Agents and Multi-Agent Systems*, pp. 251–259, 2019.
- [32] N. Tiwari, A. Ghosh, and V. Aggarwal, "Reinforcement learning for mean field game," arXiv preprint arXiv:1905.13357, 2019.
- [33] X. Guo, A. Hu, R. Xu, and J. Zhang, "Learning mean-field games," in *Neural Information Processing Systems*, pp. 4966–4976, 2019.
- [34] S. G. Subramanian, P. Poupart, M. E. Taylor, and N. Hegde, "Multi type mean field reinforcement learning," arXiv preprint arXiv:2002.02513, 2020.
- [35] M. A. uz Zaman, K. Zhang, E. Miehling, and T. Başar, "Reinforcement learning in non-stationary discrete-time linear-quadratic mean-field games," in *2020 59th IEEE Conference on Decision and Control (CDC)*, pp. 2278–2284, IEEE, 2020.
- [36] A. Angiuli, J.-P. Fouque, and M. Laurière, "Unified reinforcement Q-learning for mean field game and control problems," arXiv:2006.13912, 2020.
- [37] M. Gagrani, S. Sudhakara, A. Mahajan, A. Nayyar, and Y. Ouyang, "Thompson sampling for linear quadratic mean-field teams," arXiv preprint arXiv:2011.04686, 2020.
- [38] J. Sternby, "On consistency for the method of least squares using martingale theory," *IEEE T. on Automatic Control*, vol. 22, no. 3, pp. 346–352, 1977.