

Adaptive Online Distributed Optimal Control of Very-Large-Scale Robotic Systems

Pingping Zhu , Member, IEEE, Chang Liu , and Silvia Ferrari , Senior Member, IEEE

Abstract—Autonomous systems comprised of many cooperative agents have the potential for enabling long-duration tasks and data collection critical to the understanding of a wide range of phenomena in spatially and temporally variable environments. The adaptive distributed optimal control approach presented in this article extends online approximate dynamic programming to very-large-scale robotics (VLSR) systems that must operate and adapt to highly uncertain and variable environments. Optimal mass transport theory is used to show that, in the Wasserstein–Gaussian mixture model space, the VLSR system’s cost to go can be represented by a value functional of the robot distribution and dynamic environmental maps. The approach is demonstrated on a cooperative path planning problem in which knowledge of the obstacles in the environment changes incrementally over time based on *in situ* measurements. Numerical simulations show that the proposed approach significantly outperforms existing methods by finding an approximately optimal solution that avoids obstacles and meets a desired final robot distribution using minimum energy.

Index Terms—Adaptive, cooperative, environmental adaptation, multiagent reinforcement learning (MARL), optimal control, path planning, very-large-scale robotic (VLSR).

I. INTRODUCTION

WITH THE advent of low-cost sensors and embedded systems, very-large-scale robotic (VLSR) systems comprised of hundreds of autonomous robots are becoming a viable solution for conducting long-duration autonomous tasks over large regions of interest [1]. To date, significant progress has been made on VLSR optimal control [2]–[5]; multiagent reinforcement learning (MARL) [6]–[13]; multiagent path planning [14], [15]; and swarm robotics [16]–[21] approaches. However, these and other VLSR methods assume that knowledge of the environment is provided *a priori*. Therefore, the resulting

planning and control laws may not be adaptable to changing *in situ* conditions that autonomous robots are likely to encounter when operating over a large region for long periods of time.

Online learning and replanning are typically too computationally expensive due to the combinatorial nature of MARL [9]. In fact, even in known environments, optimal planning for N cooperative robots has been shown to be PSPACE-hard [2]. Distributed optimal control (DOC) [3]–[5] and Nash certainty equivalence (NCE) approaches [22], [23] overcome the scalability issue by defining a macroscopic state, such as the robot distribution or the robot mass, by virtue of a restriction operator and accompanying consistency relationships mapping robot kinodynamic equations onto a macroscopic evolution equation. Although the computational burden is significantly reduced [3]–[5], determining the optimal restriction operator remains too time consuming for online adaptation in response to *in situ* measurements. In VLSR long-duration applications ranging from ocean robotics to space systems, robots operate in the presence of significant uncertainties. At the same time, their performance depends on environmental conditions that cannot be accurately predicted *a priori*, thus requiring adaptation or replanning subject to online measurements [1].

A model-free MARL approach based on mean field control was recently proposed in [11] to address both scalability and uncertainty issues by trial-and-error approximations of the value function. Similarly to the DOC approach used in this article, model-free MARL employs a macroscopic state represented by the robot distribution or probability density function (PDF), and a reward (or Lagrangian) function that depends both on microscopic and macroscopic robot states. Then, the MARL approach is cast as a Markov decision process on the Wasserstein space of measures and implemented by learning a deterministic control law offline. The microscopic robot control law is given by a functional of the macroscopic and microscopic states. However, by this approach only the social average reward can be optimized and the robot decisions are myopic, namely, they are based solely on the latest robot state and robot distribution.

This article presents new adaptive DOC (ADOC) theory for VLSR systems that must operate optimally over time based on *in situ* measurements that influence immediate and future cooperative sensing and navigation performance in highly uncertain environments. The ADOC approach is assumed to be centralized for simplicity but is applicable to a distributed system of robots with a macroscopic evolution equation provided by a stochastic differential equation (SDE), also known as distributed parameter system [5]. Thus, the proposed ADOC approach can be viewed

Manuscript received January 15, 2021; revised January 17, 2021 and June 21, 2021; accepted July 1, 2021. Date of publication July 14, 2021; date of current version August 24, 2021. This work was supported in part by the Office of Naval Research Code 321, in part by the National Science Foundation under Grant ECCS-1556900 and Grant ATD-1738010, and in part by the Marshall University. Recommended by Associate Editor M. Prandini. (Corresponding author: Pingping Zhu.)

Pingping Zhu is with the Department of Computer Science and Electrical Engineering, Marshall University, Huntington, WV 25705 USA (e-mail: zhup@marshall.edu).

Chang Liu and Silvia Ferrari are with the Sibley School of Mechanical and Aerospace Engineering, Cornell University, Ithaca, NY 14853 USA (e-mail: changliu@berkeley.edu; ferrari@cornell.edu).

Digital Object Identifier 10.1109/TCNS.2021.3097306

as a reinforcement learning–approximate dynamic programming (RL-ADP) approach [24], [25] applicable to distributed parameter systems for which environmental conditions must be learned online.

The new ADOC theory presented in this article is based on several key contributions. This article shows that when the environmental information is changing over time, the system’s cost to go cannot be expressed as a value function but requires the use of a value *functional*. Then, the ADOC approach can be formulated in a Wasserstein–Gaussian mixture model (GMM) space using the optimal mass transport (OMT) theorem, by assuming that the optimal time-varying robot distribution can be learned and represented via GMMs. A fast ADOC online solution is afforded by proving that the ADOC problem amounts to solving a linear program in a subspace of the Wasserstein-GMM space. The ADOC computational complexity, convergence, and lower and upper bounds on the optimal value functional are also derived in this article. The effectiveness of ADOC is demonstrated on a VLSR obstacle avoidance problem, in which the obstacles are sensed online. The numerical results show that the ADOC approach significantly outperforms existing state-of-the-art methods and scales up to VLSR systems with hundreds of robots.

II. PROBLEM FORMULATION

Consider the problem of adaptively planning the trajectories of a VLSR system comprised of N cooperative robots deployed in a large obstacle-populated region of interest (ROI) $\mathcal{W} \subset \mathbb{R}^2$. Although the M obstacles, $\mathcal{B}_1, \dots, \mathcal{B}_M \subset \mathcal{W}$ are unknown *a priori*, an approximate map of the obstacle region $\hat{\mathcal{B}}(t)$ is provided over time by sensors and mapping algorithms on-board the N robots. At the initial time, t_0 , an obstacle map $\hat{\mathcal{B}}_0 = \hat{\mathcal{B}}(t_0)$ is provided based on prior information, and, thus, may be subject to significant errors, potentially causing blocked passages and narrow regions to be unknown *a priori*. Let the obstacle map function (OMF) be defined as a binary time-varying function, $m(\mathbf{x}, t) : \mathcal{W} \times \mathbb{R} \rightarrow \{0, 1\}$, where 1 indicates that the position $\mathbf{x} \in \mathcal{W}$ is occupied by obstacles at time t , and 0 indicates that it is unoccupied. Let \mathcal{M} denote the space of all possible OMFs in the ROI, such that $m(\cdot, t) \in \mathcal{M}$. In this article, the OMF is obtained from a Hilbert occupancy map $h(\mathbf{x}, t)$ [26], as shown in [27], such that

$$m(\mathbf{x}, t) = \begin{cases} 1, & \text{if } h(\mathbf{x}, t) > 0.5 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where $h(\cdot)$ provides the probability of obstacle occupancy.

The microscopic robot dynamics are modeled by a SDE

$$\dot{\mathbf{x}}_i(t) = \mathbf{f}[\mathbf{x}_i(t), \mathbf{u}_i(t), t] \quad (2)$$

$$\mathbf{x}_i(t_0) = \mathbf{x}_{i_0}, i = 1, \dots, N, \quad (3)$$

where $\mathbf{x}_i \in \mathcal{W}$ is the i th robot’s state or configuration, $\mathbf{u}_i(t) \in \mathcal{U}$ is the i th robot’s control input, and \mathbf{x}_{i_0} is the i th robot’s initial configuration. Also, it is assumed that the robot state is fully observable and known with negligible errors.

All robots are equipped with identical omnidirectional range sensors that allow them to update the OMF based on *in situ* measurements. The field of view (FOV) of the i th range sensor, denoted by $\mathcal{S}_i(t) \subset \mathcal{W}$, is represented by a circle of radius r centered at $\mathbf{x}_i(t)$. Then, the region covered by the sensors at time t can be represented by $\mathcal{S}(t) = \cup_{i=1}^N \mathcal{S}_i(t)$. Assuming connectivity and information sharing, obstacle presence at position $\mathbf{x} \in \mathcal{W}$ is observed and updated at time t if and only if $\mathbf{x} \in \mathcal{S}(t)$.

Let the VLSR macroscopic state be represented by a PDF, $\wp(\mathbf{x}, t) \in \mathcal{P}(\mathcal{W})$, where $\mathcal{P}(\mathcal{W})$ is the space of PDFs with support \mathcal{W} [3]–[5]. Then, the VLSR system performance over a time interval $[t_0, t_f]$ can be expressed by an integral cost function

$$J[\wp(\mathbf{x}, t)] = \phi[\wp(t_f), \wp_f] + \int_{t_0}^{t_f} \mathcal{L}[\wp(\mathbf{x}, t), m(\mathbf{x}, t)] dt \quad (4)$$

representing the cost required for the robots to move from a given initial distribution \wp_0 at time t_0 to a desired distribution \wp_f at t_f . The functionals $\phi[\wp(t_f), \wp_f]$ and $\mathcal{L}[\wp(\mathbf{x}, t), m(\mathbf{x}, t)]$ denote the terminal cost and the instantaneous cost or “Lagrangian,” respectively. Because the VLSR system’s performance depends on the changing obstacle map, $m(\cdot, \cdot)$, the Lagrangian is a functional of the OMF for all $t \in [t_0, t_f]$. The result is a new ADP problem in which the functional structure of the Lagrangian is unknown and must be learned from the range-sensor measurements over time.

III. ADAPTIVE DOC APPROACH

Because sensor measurements become available at discrete sampling times and must be fused to obtain the OMF, the DOC problem is first discretized with respect to time. Let Δt denote the time required to obtain and process the obstacle measurements, such that the time interval $[t_0, t_f]$ can be discretized into $T_f = (t_f - t_0)/\Delta t$ time steps, indexed by $t_k = t_0 + k\Delta t$, where $k = 1, \dots, T_f$. The robot dynamics in (2) can be discretized and reformulated as

$$\mathbf{x}_{k+1,i} = \mathbf{x}_{k,i} + \dot{\mathbf{x}}_{k,i} \Delta t, \quad i = 1, \dots, N \quad (5)$$

$$\dot{\mathbf{x}}_{k,i} = \mathbf{f}[\mathbf{x}_{k,i}, \mathbf{u}_{k,i}] \quad (6)$$

where $\mathbf{x}_{k,i} = \mathbf{x}_i(t_k)$ and $\mathbf{u}_{k,i} = \mathbf{u}_i(t_k)$. Then, the macroscopic cost function (4) can be approximated by

$$J(\wp_{0:T_f}) \triangleq \phi(\wp_{T_f}, \wp_f) + \sum_{k=0}^{T_f-1} \mathcal{L}(\wp_k, m_k) \quad (7)$$

where $\wp_{0:T_f} = [\wp_0 \cdots \wp_{T_f}]$, $\wp_k = \wp(\cdot, t_k)$ and $m_k = m(\cdot, t_k)$.

Because the OMF is time-varying and unknown for future times, the cost function (7) cannot be optimized with respect to the robot distribution using existing approaches [3]–[5]. An online optimization approach for this new class of finite horizon ADP problems is developed here by reformulating the objective function as follows:

$$J(\wp_{0:T}) = \phi(\wp_{T_f}, \wp_f) + \sum_{k=0}^{T_f-1} \mathcal{L}[\wp_k, m_k, \mathcal{C}(\wp_k, m_k)] \quad (8)$$

where $\mathcal{C} : \mathcal{P} \times \mathcal{M} \mapsto \mathcal{P}$ is the control law functional, and

$$\varphi_{k+1} = \mathcal{C}(\varphi_k, m_k). \quad (9)$$

In the rest of this article, $\mathcal{C}(\varphi_k, m_k)$ is abbreviated by \mathcal{C}_k , and $\mathcal{L}[\varphi_k, m_k, \mathcal{C}(\varphi_k, m_k)]$ is abbreviated by $\mathcal{L}(\varphi_k, m_k, \mathcal{C}_k)$.

Because the VLSR goal is to reach φ_f by time t_f , there exists at least one terminal macroscopic state $\varphi_T \in \mathcal{P}$ that is cost-free, and, thus, can be absorbed in the state for all $k > T$, such that $\mathcal{L}(\varphi_T, m_k, \mathcal{C}_k) = 0$ and $\mathcal{C}(\varphi_T, m_k) = \varphi_T$ when $\varphi_T = \varphi_f$, or when the distance between φ_T and φ_f is less than a user-defined threshold. Then, the VLSR control problem can be assumed to terminate upon reaching φ_T at time $T \leq T_f$.

Unlike traditional ADP problems [24], [25], the control law functional in (9) depends on the OMF, for which there exists no evolution or dynamic equation because it is updated based on exogenous measurements obtained by the robots over time. Then, it is assumed the latest map m_k is the best estimate of the obstacle layout available at time t_k . Without loss of generality, the $(1 \times T)$ vector of OMFs, $\mathbf{M}_k = [m_0 \cdots m_{k-1} m_k \cdots m_k]$, is used to represent the map history. If and when an environmental prediction model is available, the method can be easily applied by modifying the definition of \mathbf{M}_k .

A. ADOC Value Functional

The goal of the discrete-time ADOC approach is to learn the optimal control law functional, $\mathcal{C}_k^* : \mathcal{P} \mapsto \mathcal{P}$, online at every k th time step based on the OMF, m_k . Then, the optimal and nonmyopic policy associated with \mathbf{M}_k can be expressed by the vector of functionals

$$\mathbf{\Pi}_k^* \triangleq [\mathcal{C}_0^* \cdots \mathcal{C}_{k-1}^* \mathcal{C}_k^* \cdots \mathcal{C}_k^*] \quad (10)$$

and must be determined over time so as to minimize (8). By optimizing the policy subject to the robot kinodynamics (6), reachability is guaranteed and the optimal robot distribution, φ_{k+1}^* , can be realized by the robots using the microscopic control described in Section VI-A.

Given a map and policy, \mathbf{M}_k and $\mathbf{\Pi}_k$, at time step k , the ADOC value functional is defined as

$$\begin{aligned} & \mathcal{V}(\varphi_l, l | \mathbf{M}_k, \mathbf{\Pi}_k, k) \\ & \triangleq \begin{cases} \phi(\varphi_T, \varphi_f) + \sum_{\tau=l}^{k-1} \mathcal{L}(\varphi_\tau, \mathbf{M}_k(\tau), \mathcal{C}_\tau) \\ \quad + \sum_{\tau=k}^{T-1} \mathcal{L}(\varphi_\tau, \mathbf{M}_k(\tau), \mathcal{C}_k), & 0 \leq l < k \\ \phi(\varphi_T, \varphi_f) + \sum_{\tau=l}^{T-1} \mathcal{L}(\varphi_\tau, \mathbf{M}_k(\tau), \mathcal{C}_k), & k \leq l < T \\ \phi(\varphi_T, \varphi_f), & l = T \end{cases} \end{aligned} \quad (11)$$

and is abbreviated by $\mathcal{V}_k(\varphi_l, \mathbf{M}_k, \mathcal{C}_k)$ hereon. Then, from (8), it can be shown that $J(\varphi_{0:T}) = \mathcal{V}_{T-1}(\varphi_0, \mathbf{M}_{T-1})$. The ADOC cost-to-go estimate or “Q-functional” is defined as

$$\mathcal{Q}_k(\varphi_l, \mathbf{M}_k, \varphi_{l+1}) \triangleq \mathcal{L}(\varphi_l, \mathbf{M}_k(l), \varphi_{l+1}) + \mathcal{V}_k(\varphi_{l+1}, \mathbf{M}_k) \quad (12)$$

and minimized as explained in the next subsection.

B. ADOC Optimal Control Law Functional

Learning a functional operator online is computationally challenging [28], [29], and, in this article, it is approached using the critic-only Q-learning (CoQL) method developed in [30]. Applying Bellman’s equation to the ADOC value- and Q-functionals introduced in the previous subsection, the optimal Q-function can be obtained as follows:

$$\mathcal{Q}_k^*(\varphi_l, \mathbf{M}_k, \varphi_{l+1}^*) = \min_{\varphi_{l+1}} [\mathcal{L}(\varphi_l, m_k, \varphi_{l+1}) + \mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k)]$$

for any $k \leq l < T$, where $\mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k)$ is the abbreviation of the optimal value functional, $\mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k, \mathcal{C}_k^*)$. Then, the optimal robot distribution can be obtained by solving the optimization problem

$$\begin{aligned} \varphi_{k+1}^* &= \arg \min_{\varphi_{k+1}} [\mathcal{Q}_k^*(\varphi_k, \mathbf{M}_k, \varphi_{k+1})] \\ &= \arg \min_{\varphi_{k+1}} [\mathcal{L}(\varphi_k, m_k, \varphi_{k+1}) + \mathcal{V}_k^*(\varphi_{k+1}, \mathbf{M}_k)]. \end{aligned} \quad (13)$$

IV. BACKGROUND ON OPTIMAL MASS TRANSPORT

An efficient approach for the online optimization of the ADOC Q-function is developed by measuring the cost of the VLSR-distribution evolution using OMT theory [31], [32]. When compared to other approaches for measuring information divergence, such as the Kullback–Leibler (KL) divergence or the Cauchy–Schwarz (CS) divergence, OMT affords significant computational savings and also provides important metric properties [32].

Let $\varphi_1, \varphi_2 \in \mathcal{P}(\mathcal{W})$ denote two PDFs with support \mathcal{W} , and let $\Pi(\varphi_1, \varphi_2) \subset \mathcal{P}(\mathcal{W} \times \mathcal{W})$ denote the set of all joint PDFs characterized by marginals measures along the two coordinate directions that coincide with φ_1 and φ_2 , respectively, and such that

$$\begin{aligned} \Pi(\varphi_1, \varphi_2) &\triangleq \left\{ \pi \in \mathcal{P}(\mathcal{W} \times \mathcal{W}) \right. \\ &\left. \int_{\mathbf{x}_2 \in \mathcal{W}} \pi(\cdot, \mathbf{x}_2) d\mathbf{x}_2 = \varphi_1, \text{ and } \int_{\mathbf{x}_1 \in \mathcal{W}} \pi(\mathbf{x}_1, \cdot) d\mathbf{x}_1 = \varphi_2 \right\}. \end{aligned} \quad (14)$$

Then, the Wasserstein metric is defined as

$$W_2(\varphi_1, \varphi_2) \triangleq \left[\inf_{\pi \in \Pi(\varphi_1, \varphi_2)} \int_{\mathcal{W} \times \mathcal{W}} \|\mathbf{x}_1 - \mathbf{x}_2\|^2 d\pi(\mathbf{x}_1, \mathbf{x}_2) \right]^{1/2} \quad (15)$$

where $\|\cdot\|$ is the Euclidean distance [31], and can be shown finite provided the second moments of φ_1 and φ_2 exist [33].

Furthermore, if both of the marginals $\varphi_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\varphi_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ are Gaussian distributions, the Wasserstein metric can be expressed as

$$\begin{aligned} W_2(\varphi_1, \varphi_2) &= \left\{ \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2 \right. \\ &\quad \left. + \text{tr} \left[\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2 \left(\boldsymbol{\Sigma}_1^{1/2} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{1/2} \right)^{1/2} \right] \right\}^{1/2} \end{aligned} \quad (16)$$

where $\text{tr}[\cdot]$ denotes the trace of a matrix [32].

Although the Wasserstein metric of two Gaussian distributions can be obtained in closed form, there is no efficient representation for general distributions [34]. Recently in [32], a new metric in the space of all GMMs, $\mathcal{G}(\mathcal{W})$, referred to as Wasserstein-GMM (WG) metric, and defined as

$$d(\varphi_1, \varphi_2) \triangleq \left\{ \min_{\pi \in \Pi(\omega_1, \omega_2)} \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} [W_2(g_1^i, g_2^j)]^2 \pi(i, j) \right\}^{1/2} \quad (17)$$

was proposed as an efficient approximation to the Wasserstein metric for any two distributions $\varphi_1, \varphi_2 \in \mathcal{G}(\mathcal{W})$, where g_1^i and g_2^j are corresponding Gaussian mixture components with weights ω_1 and ω_2 , respectively. Then, $\Pi(\omega_1, \omega_2)$ can be used to denote the space of joint GMMs, referred to as Wasserstein-GMM space.

V. ADOC SOLUTION IN WASSERSTEIN-GMM SPACE

Hereon, let the Wasserstein metric $W_2(\varphi_k, \varphi_{k+1})$ represent the distance between two robot PDFs, φ_k and φ_{k+1} . For a fixed time interval, Δt , $W_2(\varphi_k, \varphi_{k+1})$ is proportional to the PDF velocity

$$\nu_k \triangleq \frac{W_2(\varphi_k, \varphi_{k+1})}{\Delta t} \quad (18)$$

and it can be shown that the following proportionalities hold: $E_k \propto (\nu_k)^2 \propto [W_2(\varphi_k, \varphi_{k+1})]^2$, where E_k is the energy required to move from φ_k to φ_{k+1} . Then, the WG metric in (17) can be adopted as the energy cost in the Lagrangian, provided the robot distributions obey the following assumption.

Assumption 1: Assume that within an acceptable error the optimal robot distribution can be approximated by GMMs, such that

$$\varphi_k = \sum_{i=1}^{N_k} \omega_k^i g_k^i, \quad k = 0, \dots, T \quad (19)$$

$$\varphi_f = \sum_{j=1}^L \omega_f^j g_f^j \quad (20)$$

where N_k and L are the numbers of Gaussian components g_k^i and g_f^j with means μ_k^i and μ_f^j , and covariance matrices Σ_k^i and Σ_f^j , respectively.

Considering the kernel density estimation (KDE) with Gaussian kernels as a special case of GMM [35], choose $N_k \leq N$, where a relatively small number of Gaussian components is typically required in practice to represent useful robot distributions, and, thus, $N_k \ll N$. Now, let $\mathbf{g}_k \triangleq [g_k^1 \cdots g_k^{N_k}]$ and $\mathbf{g}_f \triangleq [g_f^1 \cdots g_f^L]$ denote vectors of Gaussian components that are used to approximate the robot time-varying and final PDFs, respectively, by means of corresponding GMM weight vectors $\omega_k \triangleq [\omega_k^1 \cdots \omega_k^{N_k}]$ and $\omega_f \triangleq [\omega_f^1 \cdots \omega_f^L]$, respectively. Then, φ_k and φ_f are fully specified by the tuples of parameters $\Theta_k = (N_k, \mathbf{g}_k, \omega_k)$ and $\Theta_f = (L, \mathbf{g}_f, \omega_f)$. It also follows that

the control law functional can be expressed as

$$\mathcal{C}_{k+1} = \sum_{j=1}^{N_{k+1}} \omega_{k+1}^j g_{k+1}^j = \sum_{i=1}^{N_k} \sum_{j=1}^{N_{k+1}} \pi_k(i, j) g_{k+1}^j \quad (21)$$

where $\pi_k \in \Pi(\omega_k, \omega_{k+1})$ is the joint probability distribution and, given φ_k and m_k , the control law is fully specified by the tuple of parameters, $\Phi_k = (N_{k+1}, \mathbf{g}_{k+1}, \pi_k)$.

The VLSR performance is represented by an integral cost function of the robot PDF in the form (8), derived as follows. First, the distance between two GMM PDFs φ_k and φ_{k+1} is defined as

$$\tilde{d}(\varphi_k, m_k, \mathcal{C}_k) \triangleq \left\{ \sum_{i=1}^{N_k} \sum_{j=1}^{N_{k+1}} [W_2(g_k^i, g_{k+1}^j)]^2 \pi_k(i, j) \right\}^{1/2}$$

and is minimized to obtain the WG metric

$$d(\varphi_k, \varphi_{k+1}) = \min_{\pi_k} [\tilde{d}(\varphi_k, m_k, \mathcal{C}_k)]. \quad (22)$$

Second, knowledge of the OMF can be used for obstacles avoidance by minimizing its inner product with φ_{k+1} , such that the Lagrangian is defined as

$$\mathcal{L}(\varphi_k, m_k, \mathcal{C}_k) = [\tilde{d}(\varphi_k, m_k, \mathcal{C}_k)]^2 + \langle \mathcal{C}_k, m_k \rangle_{\mathcal{W}} \quad (23)$$

where $\langle \cdot, \cdot \rangle_{\mathcal{W}}$ is the inner product over \mathcal{W} , and $\mathcal{C}_k = \varphi_{k+1}$.

The VLSR performance at the final time depends on the robot distance to the desired PDF, φ_f , and is expressed by the terminal cost

$$\phi[\varphi(t_f), \varphi_f] = \phi(\varphi_T, \varphi_f) \triangleq [d(\varphi_T, \varphi_f)]^2. \quad (24)$$

Then, the VLSR integral cost function can be rewritten as

$$J \triangleq [d(\varphi_T, \varphi_f)]^2 + \sum_{k=0}^{T-1} \left\{ [\tilde{d}(\varphi_k, m_k, \mathcal{C}_k)]^2 + \langle \varphi_{k+1}, m_k \rangle_{\mathcal{W}} \right\}$$

and, from (11), the ADOC value functional is given by

$$\mathcal{V}_k(\varphi_k, \mathbf{M}_k, \mathcal{C}_k) = \mathcal{L}(\varphi_k, m_k, \mathcal{C}_k) + \mathcal{V}_k(\varphi_{k+1}, \mathbf{M}_k, \mathcal{C}_k). \quad (25)$$

A. ADOC Value Functional Approximation

As in classical RL-ADP approaches [24], [25], the optimal cost to go, $\mathcal{V}_k^*(\varphi_{k+1}, \mathbf{M}_k)$ in the recurrence relationship (13) is unknown and must be approximated over time. In particular, the approximation is obtained here by using an upper bound on the optimal value functional so as to derive convergence guarantees (Section VII). Consider an approximate control law, $\tilde{\mathcal{C}}_k(\varphi_\tau, m_k) : \mathcal{P} \mapsto \mathcal{P}$, obtained by holding the number of Gaussian components (N_τ) fixed, such that

$$N_{\tau+1} = N_\tau \quad (26)$$

$$\pi_\tau(i, j) = \begin{cases} \omega_\tau^i, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}, \quad i, j = 1, \dots, N_\tau \quad (27)$$

$$\varphi_{\tau+1} = \tilde{\mathcal{C}}_k(\varphi_\tau, m_k) \quad \forall \tau : k+1 \leq \tau \leq T-1. \quad (28)$$

By recursively applying the approximate control law, the evolution of the robot PDF can be generated as follows:

$$\varphi_\tau = \sum_{j=1}^{N_{k+1}} \omega_{k+1}^j g_\tau^j, \quad \tau = k+1, \dots, T \quad (29)$$

according to the N_{k+1} trajectories of the Gaussian components.

Next, consider the L Gaussian components of the desired goal PDF φ_f , denoted by the set $\{g_f^j\}_{j=1}^L$. There are $N_{k+1} \times L$ trajectories of Gaussian components that are characterized by the minimum cost function

$$\tilde{\mathcal{L}}_k^{j,j} = \begin{cases} \min \left\{ [W_2(g_T^j, g_f^j)]^2 \right. \\ \left. + \sum_{\tau=k+1}^{T-1} [W_2(g_\tau^j, g_{\tau+1}^j)]^2 \text{ if } k+1 < T \right. \\ \left. + \sum_{\tau=k+1}^{T-1} \langle g_{\tau+1}^j, m_k \rangle \right\} \\ [W_2(g_T^j, g_f^j)]^2, & \text{if } k+1 = T \end{cases} \quad (30)$$

and are indexed by j and j at the k th time step. Then, an upper bound on the optimal ADOC value functional is provided by the following theorem.

Theorem 1 (Upper bound of optimal value functional):

Given φ_{k+1} and φ_f defined in (19) and (20), respectively, there exists an upper bound of the optimal value functional $\mathcal{V}_k^*(\varphi_{k+1}, \mathbf{M}_k)$, which is denoted by $\tilde{\mathcal{V}}_k(\varphi_{k+1}, \mathbf{M}_k) = \mathcal{V}_k(\varphi_{k+1}, \mathbf{M}_k, \hat{\mathcal{C}}_k)$, such that

$$\mathcal{V}_k^*(\varphi_{k+1}, \mathbf{M}_k) \leq \tilde{\mathcal{V}}_k(\varphi_{k+1}, \mathbf{M}_k) \triangleq \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j) \quad (31)$$

where $\tilde{\pi}_k(j, j) \in \Pi(\omega_{k+1}, \omega_f)$ is the robot joint PDF.

The proof of *Theorem 1* is provided in Appendix A. Hereon, the upper bound in (31) is used to approximate the optimal value functional $\mathcal{V}_k^*(\varphi_{k+1}, \mathbf{M}_k)$.

B. Optimal ADOC Control Law

From (13), (25), and (31), the optimal control law functional at time k is given by

$$\begin{aligned} \mathcal{C}_k^* &\approx \arg \min_{\mathcal{C}_k} \left[\mathcal{L}(\varphi_k, m_k, \mathcal{C}_k) + \tilde{\mathcal{V}}_k(\varphi_{k+1}, \mathbf{M}_k) \right] \\ &= \arg \min_{\mathcal{C}_k} \left\{ [\tilde{d}(\varphi_k, m_k, \mathcal{C}_k)]^2 + \langle \mathcal{C}_k, m_k \rangle_{\mathcal{W}} \right. \\ &\quad \left. + \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j) \right\} \end{aligned} \quad (32)$$

which amounts to a calculus of variations problem. However, because the control law functional is parameterized by the tuple Φ_k , the optimal control functional can be approximated by means of the optimal parameters

$$\begin{aligned} \Phi_k^* &= \arg \min_{\Phi_k} \left\{ \sum_{i=1}^{N_k} \sum_{j=1}^{N_{k+1}} [W_2(g_k^i, g_{k+1}^j)]^2 \pi_k(i, j) \right. \\ &\quad \left. + \sum_{i=1}^{N_k} \sum_{j=1}^{N_{k+1}} \langle g_{k+1}^j, m_k \rangle_{\mathcal{W}} \pi_k(i, j) + \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j) \right\} \end{aligned} \quad (33)$$

where the upper bound in (31) is used to approximate the optimal value functional.

Imposing the following constraints on the GMM weights:

$$\omega_k^i = \sum_{j=1}^{N_{k+1}} \pi_k(i, j), \quad \omega_f^j = \sum_{j=1}^{N_{k+1}} \tilde{\pi}_k(j, j) \quad (34)$$

$$\omega_{k+1}^j = \sum_{i=1}^{N_k} \pi_k(i, j) = \sum_{j=1}^L \tilde{\pi}_k(j, j) \quad (35)$$

the optimization problem in (33) can be rewritten as

$$\Phi_k^* = \arg \min_{\Phi_k} \sum_{j=1}^{N_{k+1}} \left[\sum_{i=1}^{N_k} \mathcal{L}_k^{i,j} \pi_k(i, j) + \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j) \right] \quad (36)$$

where the energy cost associated with the motion of the Gaussian components with respect to the OMF is

$$\mathcal{L}_k^{i,j} = [W_2(g_k^i, g_{k+1}^j)]^2 + \langle g_{k+1}^j, m_k \rangle_{\mathcal{W}}. \quad (37)$$

From Φ_k^* , the optimal robot PDF, φ_{k+1}^* , is approximated by marginalizing over the approximate joint distribution $\hat{\pi}_k^*$.

Substituting $\hat{\varphi}_{k+1}^*$ in (23), the Lagrangian can be expressed as a function of the WG metric between φ_k and $\hat{\varphi}_{k+1}^*$

$$\begin{aligned} \mathcal{L}(\varphi_k, m_k, \hat{\mathcal{C}}_k^*) &= [\tilde{d}(\varphi_k, m_k, \hat{\mathcal{C}}_k^*)]^2 + \langle \hat{\mathcal{C}}_k^*(\varphi_k, m_k), m_k \rangle_{\mathcal{W}} \\ &= [d(\varphi_k, \hat{\varphi}_{k+1}^*)]^2 + \langle \hat{\varphi}_{k+1}^*, m_k \rangle_{\mathcal{W}} \end{aligned} \quad (38)$$

and the velocity of the robot PDF in (18) can be approximated using the WG metric

$$\nu_k \approx \frac{d(\varphi_k, \hat{\varphi}_{k+1}^*)}{\Delta t} \quad (39)$$

and the Lagrangian accounts for the energy consumption E_k .

VI. ADOC NUMERICAL IMPLEMENTATION

An efficient ADOC numerical solution, applicable to online adaptation by the VLSR system, is presented in this section under the following assumption.

Assumption 2: The set of Gaussian components in (19) is a subset of the union set of collocation Gaussian components $G_C = \{g_c^j\}_{j=1}^K$ and desired Gaussian components $G_f = \{g_f^j\}_{j=1}^L$

$$\{g_\tau^j\}_{j=1}^{N_\tau} \subseteq G, \quad \text{where } G \triangleq G_C \cup G_f \quad (40)$$

provided $k+1 \leq \tau \leq T$.

Based on the assumption above, \wp_τ belongs to the subspace of the GMM defined as

$$\tilde{\mathcal{G}}(\mathcal{W}, G) \triangleq \left\{ \wp \left| \begin{aligned} \wp &= \sum_{j=1}^{L+K} \omega^j g^j, \quad \sum_{j=1}^{L+K} \omega^j = 1, \\ 0 \leq \omega^j \leq 1, \quad g^j \in G, \quad j &= 1, \dots, L+K \end{aligned} \right. \right\} \quad (41)$$

where K is number of collocation components, and, thus, $N_{k+1} = L + K$. Then, the metric W_2 in (37) and (30) can be calculated in advance and, since $g_\tau^j \in G$ for all $\tau = k+1, \dots, T$, the nonlinear program (NLP) in (30) can be solved using a shortest-path algorithm on a nonnegative weighted directed graph [36], [37]. At every time step k , the directed graph, $\mathcal{G} = (G, \mathcal{E})$, is formed by assigning a node to each Gaussian component in G , and by connecting all the nodes pairwise using a set of edges \mathcal{E} . The costs associated with the edges are defined in terms of the OMF, such that

$$c_{\iota, j} = \begin{cases} [W_2(g^\iota, g^j)]^2 + \langle g^\iota, m_k \rangle_{\mathcal{W}} & \text{if } g^j \notin G_f \\ [W_2(g^\iota, g^j)]^2 & \text{if } g^j \in G_f. \end{cases} \quad (42)$$

Then, a shortest path of length of $\tilde{\mathcal{L}}_k^{j,j}$ in (30) can be found that connects node g_{k+1}^j to g_k^j .

From (36), the optimal control law functional can be approximated by solving the following optimization problem:

$$\hat{\pi}_k^* = \arg \min_{\pi_k} \sum_{j=1}^{L+K} \left[\sum_{\iota=1}^{N_k} \mathcal{L}_k^{\iota, j} \pi_k(\iota, j) + \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j, j} \tilde{\pi}_k(j, j) \right] \quad (43)$$

and the optimal robot PDF can be approximated by

$$\begin{aligned} \hat{\wp}_{k+1}^* &= \sum_{\iota=1}^{N_k} \left[\sum_{j=1}^K \hat{\pi}_k^*(\iota, j) g_c^j + \sum_{j=1}^L \hat{\pi}_k^*(\iota, j+K) g_f^j \right] \\ &= \sum_{j=1}^K (\hat{\omega}_{k+1}^j)^* g_c^j + \sum_{j=1}^L (\hat{\omega}_{k+1}^{j+K})^* g_f^j \end{aligned} \quad (44)$$

where the joint probability constraint (35) is applied. Also, for known costs of the Gaussian components motion, $\mathcal{L}_k^{\iota, j}$ and $\tilde{\mathcal{L}}_k^{j, j}$ (for all ι, j, j , and k), the approximation of the optimal value functional only depends on π_k and $\tilde{\pi}_k$ and the optimal control functional can be determined using LP algorithms.

From (44), $\hat{\wp}_{k+1}^*$ is fully specified by $L + K$ weights, many of which are equal to zero or small in magnitude. Therefore, the computational complexity can be significantly reduced by neglecting Gaussian components with weights below a user-specified threshold, and only the remaining components are normalized to approximate the optimal robot PDF.

The computational complexity of the algorithm for approximating the ADOC control law (analyzed in Section VII and shown in Table II) can be further reduced for online implementation by adopting the following assumption.

Assumption 3: The Gaussian components used to approximate the control law $\wp_{\tau+1} = \mathcal{C}_k(\wp_\tau)$, at every time step $\tau =$

TABLE I
PERFORMANCE COMPARISON

Approach	T	Solution Time (min)	$\bar{D}(0)$ (km)	$\bar{E}(T)$ (J/kg)
ADOC	701	37.0544	28.5623	545.1525
PDF-APF	2000	60.2447	99.95	1928.0478
SAPF	2000	19.6193	100.8671	3747.0543
SPP	1414	141.2573	41.9044	808.1635

TABLE II
COMPUTATIONAL COMPLEXITY COMPARISON

Algorithm		Microscopic Control
ADOC		$O(N \cdot N_k + N \cdot B_k + N^2)$
PDF-APF		$O(N \cdot L + N \cdot B_k + N^2)$
SAPF		$O(N \cdot B_k + N^2)$
SPP		$O(N \cdot B_k + N^2)$
Algorithm		Path Planning
ADOC		$O(\mathcal{E} + G \log G) N_k \cdot \mathcal{I}_k \cdot L$
SPP		$O(\mathcal{E}_i^{SPP} + G_i^{SPP} \log G_i^{SPP}) N$

$k, \dots, T-1$, are characterized by a transportation distance below a user-defined positive threshold, d_{th} , such that

$$\mathcal{E} = \{e_{\iota, j} | g^\iota, g^j \in G \text{ and } W_2(g^\iota, g^j) < d_{th}\} \quad (45)$$

and

$$\begin{aligned} \pi_\tau(\iota, j) &= 0 \text{ if } W_2(g_\tau^\iota, g_{\tau+1}^j) > d_{th} \\ \iota &= 1, \dots, N_\tau, \text{ and } j = 1, \dots, L+K \end{aligned} \quad (46)$$

where the joint probability π_τ is an $N_\tau \times (L+K)$ matrix, and g_τ^ι and $g_{\tau+1}^j$ are the Gaussian components approximating the robot PDF at consecutive times τ and $(\tau+1)$, respectively.

Now, let the index set of the Gaussian components used in the control law approximation be denoted by

$$\mathcal{I}_\tau^\iota = \{j | j \in \mathcal{I}, g^j \in G, \text{ and } W_2(g_\tau^\iota, g^j) \leq d_{th}\} \subset \mathcal{I} \quad (47)$$

where $\mathcal{I} = \{1, \dots, L+K\}$ is the index set of all components in G . Then, from Assumption 3, (43) can be rewritten as

$$\hat{\pi}_k^* = \arg \min_{\pi_k} \sum_{\iota=1}^{N_k} \sum_{j=1}^L \left\{ \sum_{j \in \mathcal{I}_k^\iota} \left[\mathcal{L}_k^{\iota, j} \pi_k(\iota, j) + \tilde{\mathcal{L}}_k^{j, j} \tilde{\pi}_k(j, j) \right] \right\} \quad (48)$$

where

$$\pi_k(\iota, j) = 0, \quad \iota = 1, \dots, N_k \text{ and } j \notin \mathcal{I}_k^\iota \quad (49)$$

$$\tilde{\pi}_k(j, j) = 0, \quad j = 1, \dots, L \text{ and } j \in \mathcal{I}_k^C \quad (50)$$

and

$$\mathcal{I}_k^C = \{j | j \in \mathcal{I} \text{ and } j \notin \mathcal{I}_k\}, \quad \mathcal{I}_k \triangleq \cup_{\iota=1}^{N_k} \mathcal{I}_k^\iota \quad (51)$$

such that (34)–(35) are satisfied. As a result, the LP in (48) can be solved using only $|\mathcal{I}_k| \times L$ shortest-paths, where “ $|\cdot|$ ” is the cardinality operator.

Because the robot PDF velocity (39) depends on the collocation Gaussian components, the distance between two sequential robot distributions is divided evenly by a user-defined interval,

D , and other robot distributions are obtained via interpolation. In particular, let the PDF distance $d(\wp_k, \hat{\wp}_{k+T_k}^*)$ be divided into $T_k = \lceil d(\wp_k, \hat{\wp}_{k+T_k}^*)/D \rceil$ intervals, where “ $\lceil \cdot \rceil$ ” is the ceiling operator. Then, $(T_k - 1)$ robot PDFs can be obtained via interpolation between each pair of Gaussian components, g_k^i and $g_{k+T_k}^j$, $i = 1, \dots, N_k$ and $j = 1, \dots, N_{k+T_k}$, as shown in [32], and the robot PDF velocity can be approximated by

$$\nu_\tau \approx \frac{d(\wp_k, \hat{\wp}_{k+T_k}^*)}{T_k \cdot \Delta t} \approx \frac{D}{\Delta t} \triangleq \bar{v}, \quad k \leq \tau \leq k + T_k. \quad (52)$$

Then, the robot PDFs travel with a relatively smooth velocity field and when $D \ll d(\wp_k, \hat{\wp}_{k+T_k}^*)$ the PDF velocity can be treated as a constant.

Furthermore, considering that there are T_k time steps from \wp_k to $\hat{\wp}_{k+T_k}^*$, the costs associated with the edges in (42) is modified as

$$c_{i,j} = \begin{cases} [W_2(g^i, g^j)]^2 + \langle g^j, m_k \rangle_{\mathcal{W}} & \text{if } g^j \notin G_f \\ + \sum_{n=1}^{T_k-1} \langle g_{i,j}^n, m_k \rangle_{\mathcal{W}} & \\ [W_2(g^i, g^j)]^2 & \text{if } g^j \in G_f \end{cases}$$

where $g_{i,j}^n$, $n = 1, \dots, T_k - 1$, are obtained via interpolation between g^i and $g^j \in G$.

A. Optimal Robot Control Law

Once the optimal evolution of the VLSR macroscopic state, represented here by the time-varying robot PDF, is obtained at a time step k , it can be used by each robot to compute a local optimal control law. For cooperative robots, the optimal PDF depends on the robot relative positions, and, thus, each robot control law also depends on the positions of other robots in the system. A centralized artificial potential field (APF) approach is adopted here that computes the microscopic robot control inputs for all N robots, represented by the vector $\mathbf{U}_k = [\mathbf{u}_{k,1}^T \cdots \mathbf{u}_{k,N}^T]^T$, so as to meet the desired PDF $\hat{\wp}_{k+1}^*$ at the next time step, based on the observed robot microscopic states $\mathbf{X}_k = [\mathbf{x}_{k,1}^T \cdots \mathbf{x}_{k,N}^T]^T$.

The attractive potential designed to “push” the robots toward the desired PDF is given by

$$U_{\text{att}} = \int_{\mathcal{W}} [\hat{\wp}_{k+1}^*(\mathbf{x}) - \gamma \tilde{\wp}_{k+1}(\mathbf{x}; \mathbf{X}_k, \mathbf{U}_k)]^2 d\mathbf{x} \quad (53)$$

where $\tilde{\wp}_{k+1}(\cdot)$ is the predicted robot PDF at time $(k+1)$ conditioned on the observed robots’ positions and controls (\mathbf{X}_k and \mathbf{U}_k). γ is a scalar parameter that determines the scattering strength of the robots chosen, such that $0 < \gamma \leq 1$. The KDE method can be used to estimate the robot PDF at the present time (k) based on the VLSR observed position vector \mathbf{X}_k [4], [5].

The repulsive potential is designed to “pull” the robots away from obstacles and from each other, in order to avoid collisions. Let ρ_i denote the minimum Euclidian distance between the robot position $\mathbf{x}_{k,i}$ and the obstacle region $\hat{B}(k)$ estimated from the latest OMF, m_k . Then, the repulsive obstacle potential for the i th robot is

$$U_{\text{obs}}^i = \frac{1}{2} \left(\frac{1}{\rho_i} - \frac{1}{\varrho} \right)^2 \cdot 1(\rho_i, \varrho)$$

where ϱ is a distance threshold used to create a region of influence within which obstacles repel robots, and $1(\rho_i, \varrho)$ is an indicator function that equals one if $\rho_i \leq \varrho$, and equals zero otherwise. Also, let $\rho_{i,\ell} = \|\mathbf{x}_{k+1,i} - \mathbf{x}_{k+1,\ell}\|$ represent the distance between predicted robot positions $\mathbf{x}_{k+1,i}$ and $\mathbf{x}_{k+1,\ell}$ ($i \neq \ell$). Then, a repulsive potential between robots is obtained as follows:

$$U_{\text{rob}}^{i,\ell} = \frac{1}{2} \left(\frac{1}{\rho_{i,\ell}} - \frac{1}{\varphi} \right)^2 \cdot 1(\rho_{i,\ell}, \varphi) \quad (54)$$

where φ is a distance threshold used to create a region of influence within which each robot repels other robots.

Then, the total VLSR potential field can be obtained from a weighted combination of the attractive and repulsive potentials

$$U = w_1 \cdot U_{\text{att}} + w_2 \cdot U_{\text{rep}} \quad (55)$$

where w_1 and w_2 are user-defined weights representing the desired tradeoff between attractive and repulsive objectives, and

$$U_{\text{rep}} = \sum_{i=1}^N U_{\text{obs}}^i + \sum_{1 \leq i \neq \ell \leq N} U_{\text{rob}}^{i,\ell}. \quad (56)$$

Once the potential field in (55) is determined from the desired robot PDF, the robot control inputs at time k are computed according to the (microscopic) control law

$$\mathbf{u}_{k,i} = -\frac{\partial U}{\partial \mathbf{u}_{k,i}}, \quad i = 1, \dots, N \quad (57)$$

which is designed to minimize the potential function U by a local gradient-based approach, summarized in Algorithm 1.

VII. ADOC ALGORITHM ANALYSIS

The ADOC approach is designed based on a novel value functional defined in (11), where the policy is updated at each time step, indexed by k . Because the ADOC solution is obtained online, the optimal policy $\mathbf{\Pi}_k^*$ is obtained incrementally by minimizing the value functional in (11) with respect to \mathcal{C}_k , where the control law history $\mathbf{\Pi}_{k-1}^*$ is given. This observation, along with the results in the previous section, is used to find a lower bound on the optimal value functional according to the following theorem.

Theorem 2 (Lower bound on optimal value functional):

Given the OMF \mathbf{M}_k at the k th time step, the optimal value functional $\mathcal{V}_k^*(\wp_l, \mathbf{M}_k)$ provides a lower bound for all optimal value functionals obtained during previous time steps, such that

$$\mathcal{V}_k^*(\wp_l, \mathbf{M}_k) \leq \mathcal{V}_q^*(\wp_l, \mathbf{M}_k), \quad 0 \leq q \leq k \text{ and } 0 \leq l \leq T$$

for all $\wp_l \in \mathcal{P}(\mathcal{W})$.

The proof of Theorem 2 is provided in Appendix B.

Now, let

$$\tilde{J}(\mathbf{\Pi}_k) \triangleq \mathcal{V}_k(\wp_0, \mathbf{M}_{T-1}), \quad 0 \leq k < T \quad (58)$$

denote the ADOC cost function associated with $\mathbf{\Pi}_k$ given \wp_0 and \mathbf{M}_{T-1} . From (8), it can be shown that

$$J(\wp_{0:T}) = \tilde{J}(\mathbf{\Pi}_{T-1}) = \mathcal{V}_{T-1}(\wp_0, \mathbf{M}_{T-1}) \quad (59)$$

Algorithm 1: ADOC Numerical Implementation.**Initialization:**

Initialize the set of K collocation Gaussian components,

G_C

Construct the directed graph $\mathcal{G} = (G, \mathcal{E})$

Let $T_k = 0$, $l = 0$, and $k = 0$

Procedure:

- 1: **while** ($k \leq T$) **and** ($\wp_k \neq \wp_f$) **do**
- 2: Update m_k using environmental observations
- 3: **if** $l == T_k$ **then**
- 4: Compute $\cup_{i=1}^{N_k} \{\mathcal{L}_k^{i,j} | j \in \mathcal{I}_k\}$ from (37)
- 5: Update \mathcal{G} according to (53)
- 6: Compute $\{\tilde{\mathcal{L}}_k^{j,j} | j \in \mathcal{I}_k \text{ and } j = 1, \dots, L\}$ by finding shortest path in \mathcal{G}
- 7: Obtain $\hat{\pi}_k^*$ by solving the LP problem in (48)
- 8: Approximate $\hat{\wp}_{k+T_k}^*$ from $\hat{\pi}_k^*$ according to (44)
- 9: Let $T_k = \lceil d(\wp_k, \hat{\wp}_{k+T_k}^*)/D \rceil$
- 10: $l = 0$
- 11: Generate robot PDFs by interpolating between \wp_k and $\hat{\wp}_{k+T_k}^*$
- 12: **end if**
- 13: Update the interpolation offset step, $l = l + 1$
- 14: Update the next robot PDF, $\hat{\wp}_{k+1}^*$
- 15: Generate \mathbf{U}_k according to (57)
- 16: Update \mathbf{X}_{k+1} according to (5) and (6)
- 17: Update time step, $k = k + 1$
- 18: **end while**

and that the ADOC cost function associated with $\mathbf{\Pi}_k^*$ for $0 \leq k < T$ converges to an upper bound of the minimum cost function $J(\wp_{0:T}^*)$, as stated in the following corollary.

Corollary 2.1: For any two ADOC optimal policies, $\mathbf{\Pi}_k^*$ and $\mathbf{\Pi}_q^*$, where $0 \leq q \leq k \leq (T - 1)$, the following inequalities hold:

$$J(\wp_{0:T}^*) \leq \tilde{J}(\mathbf{\Pi}_{T-1}^*) \leq \tilde{J}(\mathbf{\Pi}_k^*) \leq \tilde{J}(\mathbf{\Pi}_q^*) \quad (60)$$

and, thus, it follows that the optimal ADOC cost function $\tilde{J}(\mathbf{\Pi}_k^*)$ converges monotonically to $\tilde{J}(\mathbf{\Pi}_{T-1}^*)$, which is an upper bound of the optimal cost function.

Corollary 2.1 can be easily proven from (58) and Theorem 2, by applying the optimal control law functional in (13). $J(\wp_{0:T}^*)$ is obtained by minimizing (8) with respect to $\wp_{0:T}$, while the policy $\mathbf{\Pi}_k^*$ is improved incrementally by updating only one control law functional at every time step, up to time $(T - 1)$. Hence, $\mathbf{\Pi}_{T-1}^*$ is a suboptimal solution.

Furthermore, the ADOC approach relies on approximating the upper bound of the value functional (Theorem 1). The optimal control law is approximated at each time step by solving (36). To reduce the computational complexity and obtain a tractable solution, the upper bound of the value functional is calculated only once per OMF, and is not updated iteratively as in conventional ADP approaches. Finally, the performance of the ADOC approach also depends on the choice of collocation Gaussian components. Although the uniform grid adopted in this

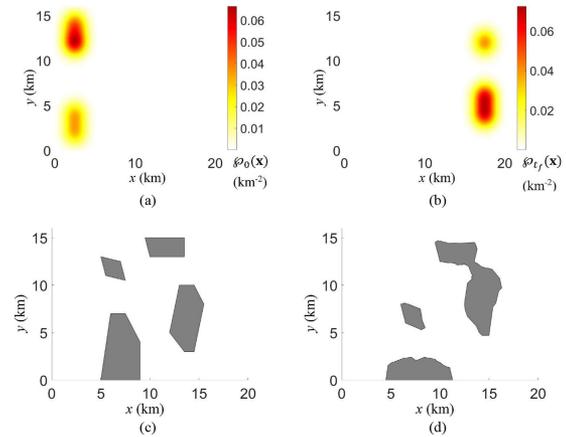


Fig. 1. VLSR system must travel from the initial robot distribution in (a) to the goal robot distribution in (b) avoiding obstacles sensed *in situ*, starting with the initial OMF (c) and learning the actual OMF (d) online.

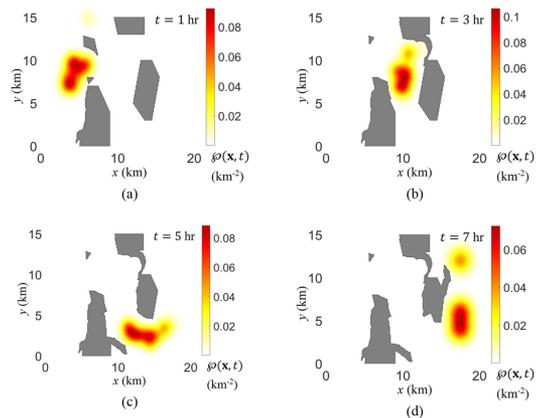


Fig. 2. Evolution of robot PDF optimized by ADOC, shown at four sample moments in time, where the gray regions represent the obstacle map updated online based on onboard sensor measurements.

article provides good performance, adaptive and multiresolution methods will be the subject of future research.

VIII. SIMULATIONS AND RESULTS

The effectiveness of the ADOC approach is demonstrated on a VLSR system comprised of $N = 500$ mobile robots characterized by single-integrator dynamics

$$\dot{\mathbf{x}}_i(t) = \mathbf{u}_i(t), \quad \mathbf{x}_i(t_0) = \mathbf{x}_{i_0}, \quad i = 1, \dots, N \quad (61)$$

where $\mathbf{x}_i = [x_i \ y_i]^T$ is the robot position, x_i and y_i are the robot xy -coordinates in inertial frame, and the control input \mathbf{u}_i is a vector of linear velocities in the x - and y -directions. The VLSR system must travel from a given initial distribution $\wp_0 = \sum_{i=1}^Q \omega_i^0 g_i^0$ to a desired distribution $\wp_f = \sum_{j=1}^L \omega_j^f g_j^f$ while avoiding collisions with partially unknown or uncertain obstacles in the ROI, $\mathcal{W} = [0, W] \times [0, H]$, where $Q = 4$, $L = 3$, $W = 20$ km, and $H = 16$ km.

The initial and desired robot PDFs are shown in Fig. 1(a) and (b), respectively. At the initial time t_0 , 500 robots characterized

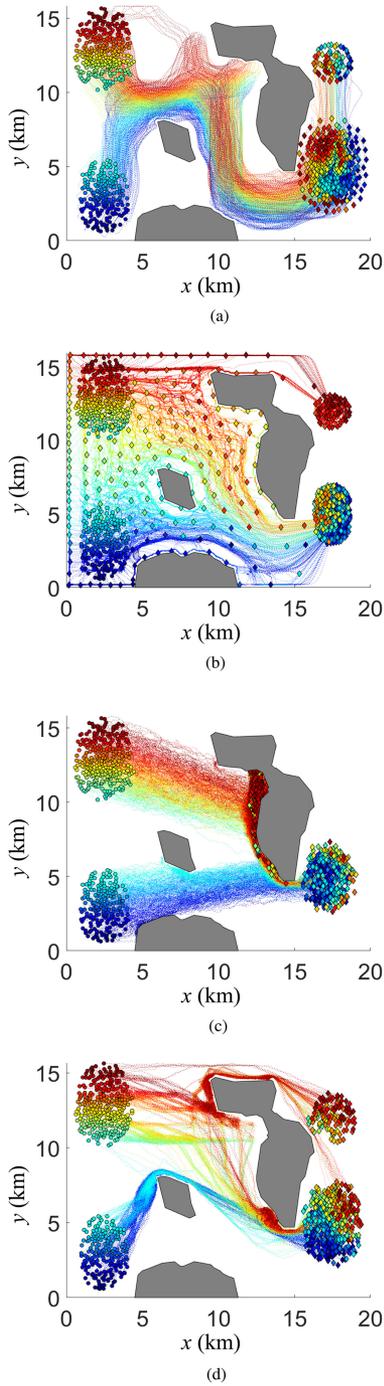


Fig. 3. Robot trajectories and final positions (diamonds) obtained by (a) ADOC, (b) PDF-APF, (c) SAPF, and (d) SPP, for the same set of initial positions (circles), and the same time-varying OMF (grey region).

by the distribution φ_0 [Fig. 3(a)] begin traveling through the ROI based on the OMF shown in Fig. 1(c). The actual obstacle region, shown in Fig. 1(d) is not yet available to them and, for example, the robots are unaware of the blocked passage in the upper-right corner of the ROI. Using an onboard range sensor characterized by an FOV with $r = 1$ km, the robots are able to update their knowledge of the obstacles and fuse their measurements to update the OMF incrementally over time.

An online ADOC solution to this VLSR path planning problem is obtained by using the set of K collocation Gaussian components

$$G_c = \left\{ g_c = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mid \boldsymbol{\mu} = [\xi - 0.5, \zeta - 0.5] \text{ km} \right.$$

$$\text{for } \xi = 1, 2, \dots, W \text{ and } \zeta = 1, 2, \dots, H$$

$$\left. \text{and } \boldsymbol{\Sigma} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \text{ km}^2 \right\}. \quad (62)$$

The component distance threshold is chosen as $d_{th} = 4$ km, and time is discretized using an interval $\Delta t = 0.01$ hr. A relatively constant distribution velocity in (52) is obtained by letting $D = 0.05$ km. Finally, the individual (microscopic) control law is obtained by constructing attractive and repulsive potentials with user-defined parameters $\gamma = 0.85$, $\varrho = 0.3$ km, and $\varphi = 0.1$ km.

An example of ADOC solution is plotted in Fig. 2 along with the OMF, which is updated based on *in situ* measurements at every time step Δt and is shown by the grey obstacle region. It can be seen in Fig. 2(b) that, after approximately three hours, the blocked passage is observed by robots. Thanks to the ADOC approach, the VLSR system is able to adapt in real time and find a new optimal solution that allows the robots to reach the desired distribution efficiently [Fig. 2(c)–(d)].

A. Performance Comparison

To the best of the authors' knowledge none of the existing VLSR methods are applicable to the online planning and control problem tackled in this article (described in Section II). For comparison, three methods are developed by extending state-of-the-art techniques referred to as PDF-based APF (PDF-APF), sampling-based APF (SAPF), and sampling-based path-planning (SPP). In the PDF-APF approach, an attractive potential field is generated by replacing $\hat{\varphi}_{k+1}^*(\mathbf{x})$ in (53) with the desired PDF φ_f . Subsequently, APF robot control laws are obtained by means of the gradient-descent method shown in (57), based on the same online OMF used by the ADOC algorithm.

In the SAPF approach, the desired final robot positions, denoted by the set $\mathcal{X}_f = \{\mathbf{x}_{f,i}\}_{i=1}^N$, are first obtained by sampling φ_f . Subsequently, since the robots are interchangeable, the desired positions are used to generate individual robot attractive potentials, and, thus, to control the robots independently of each other. In the SPP approach, the robot positions in the set $\mathcal{X}_f = \{\mathbf{x}_{f,i}\}_{i=1}^N$ are paired with the initial robot positions in $\mathcal{X}_0 = \{\mathbf{x}_i(t_0)\}_{i=1}^N$ based on the shortest relative distance. Subsequently, a shortest-path algorithm is used at every time step (k) to find the best robot trajectory based on the latest OMF, m_k . All methods avoid collisions by means of the repulsive potentials in (56). For comparison, the robot velocities are chosen to abide to $\|\mathbf{u}_i\| = 5$ km/hr for all i , although in principle they too could be optimized via ADOC in order to minimize energy consumption. Similarly, the same user-defined parameters, ϱ , φ , and d_{th} , are used in all four methods and the same maximum time $T_f = 2000$ is adopted.

Fig. 3(a)–(d) shows the VLSR trajectories obtained for the 500 robots using ADOC, PDF-APF, SAPF, and SPP. It can be

seen that using ADOC and SPP the robots are able to reach the goal distribution in the allotted time, while with PDF-APF and SAPF many of the robots fall behind or remain stuck nearby some of the obstacles found *in situ*. Table I shows the VLSR performance obtained by the four methods. All of the simulations are conducted on the same computer with an 18-core CPU and 16 G RAM. It can be seen that the solution time, T , required by the robots to reach the goal PDF is significantly lower for ADOC than for all other methods. The total runtime is also significantly reduced by ADOC despite its ability to optimize the VLSR system performance (not afforded by other methods). The optimized ADOC performance can be assessed by computing the average cost-to-go

$$\bar{D}(k) = \frac{1}{N} \sum_{i=1}^N \sum_{\tau=k}^{T-1} \|\mathbf{x}_i[(\tau+1)\Delta t] - \mathbf{x}_i(\tau\Delta t)\| \quad (63)$$

and the average energy-cost per kg

$$\bar{E}(k) = \frac{\eta}{2N} \sum_{i=1}^N \sum_{\tau=1}^k \left[\frac{\|\mathbf{x}_i(\tau\Delta t) - \mathbf{x}_i[(\tau-1)\Delta t]\|}{\Delta t} \right]^2 \quad (64)$$

where η is a unit-conversion factor. As shown in Table I, despite using the same knowledge of the obstacle region, by leveraging new optimality conditions derived from Bellman's optimality principle, the ADOC approach allows the cooperative robots to minimize time, distance traveled, and energy consumption.

B. Computational Complexity Comparison

The results shown in Table I show that the time required to obtain the ADOC solution is significantly lower than the time required by other methods of solution. This section derives the computational complexity of each algorithm, summarized in Table II. As a first step, the computational complexity required by the construction of the artificial potentials at the k th time step is analyzed. Let B_k denote the number of collocation points evenly sampled on the updated OMF contours at the k th time step, as well as on the ROI boundaries, $\partial\mathcal{W}$. Then, the (microscopic) robot control law computation at the k th time step requires the times shown in Table II, based on both repulsive and attractive potential fields. Because $N_k \ll N$ and $L \ll N$, computing the robot control law requires time $O(N^2)$ for all four methods. The square power results from the artificial potential calculation and, therefore, when robots are outside the potential's region of influence the computational complexity decreases to $O(N)$.

From [36] and [37], the computational complexity of searching a directed graph $\mathcal{G} = (G, \mathcal{E})$ for the shortest path connecting two of its nodes is $O(|\mathcal{E}| + |G| \log |G|)$. Letting $\mathcal{E}_i^{\text{SPP}}$ and G_i^{SPP} denote the set of arcs and nodes generated by the SPP algorithm for every robot i , the computation required by the path-planning ADOC and SPP methods at every time step k is found to be as shown in Table II. Therefore, when $N_k \cdot |\mathcal{I}_k| \cdot L < N$, ADOC requires significantly less computation than SPP, because typically $L \ll K$. Furthermore, since N_k , $|\mathcal{I}_k|$, and L are independent of the number of robots (N), the ADOC approach

developed in this article is scalable to very large systems of cooperative robots.

IX. CONCLUSION

This article develops a novel adaptive optimal control approach, referred to as ADOC, that is applicable to online cooperative VLSR systems applications, such as sensing and navigation. Unlike existing adaptive dynamic programming and adaptive control approaches, ADOC is developed to solve environmental adaptation problems in which the system performance, represented by the Lagrangian of the cost function, changes over time due to *in situ* measurements and observations. By optimizing the spatio-temporal evolution of the VLSR macroscopic state, such as the robot PDF, the ADOC approach provides online solutions that scale up to very large numbers of cooperative robots. The novel technical contributions in this article show that the online adaptive control of multiscale dynamical systems can be formulated as a new adaptive dynamic programming problem in the Wasserstein-GMM space, thus allowing for the application of OMT theory. The numerical simulation results presented in this article show that ADOC not only outperforms other VLSR planning methods, by optimizing the macroscopic system performance incrementally over time but it also reduces the solution time and the total amount of time required by the robots to complete the desired task.

APPENDIX A

UPPER BOUND OF OPTIMAL VALUE FUNCTIONAL

Proof: First, consider the case of $k+1 < T$. According to (25), the value functional $\mathcal{V}_k(\wp_{k+1}, m_k, \tilde{\mathcal{C}}_k)$ associated to $\tilde{\mathcal{C}}_k$, which is described in (26) and (27), can be expressed by

$$\begin{aligned} \mathcal{V}_k(\wp_{k+1}, \mathbf{M}_k, \tilde{\mathcal{C}}_k) &= [d(\wp_T, \wp_f)]^2 \\ &+ \sum_{\tau=k+1}^{T-1} \left[\tilde{d}(\wp_\tau, m_k, \tilde{\mathcal{C}}_k) \right]^2 + \sum_{\tau=k+1}^{T-1} \langle \wp_{\tau+1}, m_k \rangle_{\mathcal{W}}. \end{aligned} \quad (65)$$

According to (22), the following inequality is obtained:

$$[d(\wp_T, \wp_f)]^2 \leq \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L [W_2(g_T^j, g^j)]^2 \tilde{\pi}_k(j, j). \quad (66)$$

By recursively applying (26) and (27), the term, $[\tilde{d}(\wp_\tau, m_k, \tilde{\mathcal{C}}_k)]^2$, $k+1 \leq \tau \leq T-1$, can be expressed as

$$\begin{aligned} \left[\tilde{d}(\wp_\tau, m_k, \tilde{\mathcal{C}}_k) \right]^2 &= \sum_{j=1}^{N_{k+1}} \sum_{\iota=1}^{N_{k+1}} [W_2(g_\tau^j, g_{\tau+1}^\iota)]^2 \pi_k(j, \iota) \\ &= \sum_{j=1}^{N_{k+1}} [W_2(g_\tau^j, g_{\tau+1}^j)]^2 \omega_{k+1}^j \\ &= \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L [W_2(g_\tau^j, g_{\tau+1}^j)]^2 \tilde{\pi}_k(j, j). \end{aligned} \quad (67)$$

Substituting (66) and (67) into (65), one can have

$$\begin{aligned} \mathcal{V}_k(\varphi_{k+1}, \mathbf{M}_k, \tilde{\mathcal{C}}_k) &\leq \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \left\{ [W_2(g_T^j, g^j)]^2 \right. \\ &+ \left. \sum_{\tau=k+1}^{T-1} [W_2(g_\tau^j, g_{\tau+1}^j)]^2 + \sum_{\tau=k+1}^{T-1} \langle g_\tau^j, m_k \rangle_{\mathcal{W}} \right\} \tilde{\pi}_k(j, j). \end{aligned} \quad (68)$$

Because (68) holds for any trajectories of Gaussian components from g_{k+1}^j to g^j , $j = 1, \dots, N_{k+1}$ and $j = 1, \dots, L$, the following inequality can be obtained:

$$\begin{aligned} \mathcal{V}_k(\varphi_{k+1}, \mathbf{M}_k, \tilde{\mathcal{C}}_k) &\leq \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \min \left\{ [W_2(g_T^j, g^j)]^2 \right. \\ &+ \left. \sum_{\tau=k+1}^{T-1} [W_2(g_\tau^j, g_{\tau+1}^j)]^2 + \sum_{\tau=k+1}^{T-1} \langle g_\tau^j, m_k \rangle_{\mathcal{W}} \right\} \tilde{\pi}_k(j, j) \\ &= \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j). \end{aligned} \quad (69)$$

Next, consider the case of $k+1 = T$. According to (66), the upper bound of $\mathcal{V}_k(\varphi_T, \mathbf{M}_k, \tilde{\mathcal{C}}_k)$ is expressed by

$$\begin{aligned} \mathcal{V}_k(\varphi_T, \mathbf{M}_k, \tilde{\mathcal{C}}_k) &\leq \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L [W_2(g_T^j, g^j)]^2 \tilde{\pi}_k(j, j) \\ &= \sum_{j=1}^{N_{k+1}} \sum_{j=1}^L \tilde{\mathcal{L}}_k^{j,j} \tilde{\pi}_k(j, j). \end{aligned} \quad (70)$$

Finally, considering the definition of the optimal value functional, the theorem is proved.

APPENDIX B

LOWER BOUND OF OPTIMAL VALUE FUNCTIONAL

Proof: First, consider the case of $q = k - 1$, $0 < k < T$, and $k \leq l < T$. From any robot PDF φ_l , by using the optimal control law functional obtained at the q th time step, \mathcal{C}_q^* , recursively, a trajectory of robot PDFs $\{\varphi_\tau\}_{\tau=l}^T$ are generated associated with the OMF m_k .

In addition, consider the optimal value functional $\mathcal{V}_k^*(\varphi_l, \mathbf{M}_k)$, $k \leq l < T$. According to (11) and the Bellman equation, $\mathcal{V}_k^*(\varphi_l, \mathbf{M}_k)$ can be expressed by

$$\begin{aligned} \mathcal{V}_k^*(\varphi_l, \mathbf{M}_k) &= \min_{\varphi_{l+1}} [\mathcal{L}(\varphi_l, m_k, \varphi_{l+1}) + \mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k)] \\ &= \mathcal{L}(\varphi_l, m_k, \varphi_{l+1}^*) + \mathcal{V}_k^*(\varphi_{l+1}^*, \mathbf{M}_k) \\ &\leq \mathcal{L}(\varphi_l, m_k, \varphi_{l+1}) + \mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k) \end{aligned} \quad (71)$$

where $\varphi_{l+1} = \mathcal{C}_q^*(\varphi_l, m_k)$ and $\varphi_{l+1}^* = \mathcal{C}_k^*(\varphi_l, m_k)$.

By recursively utilizing (71), the following inequality is obtained:

$$\begin{aligned} \mathcal{V}_k^*(\varphi_l, \mathbf{M}_k) &\leq \mathcal{L}(\varphi_l, m_k, \varphi_{l+1}) + \mathcal{V}_k^*(\varphi_{l+1}, \mathbf{M}_k) \\ &\leq \mathcal{L}(\varphi_l, m_k, \varphi_{l+1}) \end{aligned}$$

$$\begin{aligned} &+ \mathcal{L}(\varphi_{l+1}, m_k, \varphi_{l+2}) + \mathcal{V}_k^*(\varphi_{l+2}, \mathbf{M}_k) \\ &\dots \\ &\leq \sum_{\tau=l}^{T-1} \mathcal{L}(\varphi_\tau, m_k, \varphi_{\tau+1}) + \mathcal{V}_k^*(\varphi_T, \mathbf{M}_k) \\ &= \sum_{\tau=l}^{T-1} \mathcal{L}(\varphi_\tau, m_k, \varphi_{\tau+1}) + \mathcal{V}_q^*(\varphi_T, \mathbf{M}_k) \\ &= \mathcal{V}_q^*(\varphi_l, \mathbf{M}_k), \quad k \leq l < T \end{aligned} \quad (72)$$

where $\mathcal{V}_k^*(\varphi_T, \mathbf{M}_k) = \mathcal{V}_q^*(\varphi_T, \mathbf{M}_k) = [d(\varphi_T, \varphi_f)]^2$ according to (11). Because this holds for $l = T$ as well, according to (11), (72) can be rewritten by

$$\mathcal{V}_k^*(\varphi_l, \mathbf{M}_k) \leq \mathcal{V}_q^*(\varphi_l, \mathbf{M}_k), \quad k \leq l \leq T. \quad (73)$$

Next, consider the case of $q = k - 1$ and $0 \leq l \leq q$. Assume that the optimal policy Π_q exists, which includes a sequence of optimal control law functionals \mathcal{C}_τ^* , $\tau = l, \dots, q$. By using these optimal control law functionals sequentially, a trajectory of robot PDFs $\{\varphi_\tau\}_{\tau=l}^q$ is generated associated with $\{m_\tau\}_{\tau=l}^q$. Moreover, similarly, by using the optimal control law functional \mathcal{C}_q^* recursively, a trajectory of robot PDFs $\{\varphi_\tau\}_{\tau=q}^T$ is generated from φ_q to φ_T associate with m_k . Thus, a trajectory of robot PDFs, $\{\varphi_\tau\}_{\tau=l}^T$ is generated from φ_l to φ_T .

Again, according to (11) and the Bellman equation, $\mathcal{V}_k^*(\varphi_l, \mathbf{M}_k)$ can be expressed by

$$\begin{aligned} \mathcal{V}_k^*(\varphi_l, \mathbf{M}_k) &= \sum_{\tau=l}^q \mathcal{L}(\varphi_\tau, m_\tau, \mathcal{C}_\tau^*) + \mathcal{V}_k^*(\varphi_k, \mathbf{M}_k) \\ &\leq \sum_{\tau=l}^q \mathcal{L}(\varphi_\tau, m_\tau, \mathcal{C}_\tau^*) + \mathcal{V}_q^*(\varphi_k, \mathbf{M}_k) \\ &= \mathcal{V}_q^*(\varphi_l, \mathbf{M}_k), \quad 0 \leq l \leq q \end{aligned} \quad (74)$$

where the inequality is obtained by applying (73). Merging (73) and (74), it is shown that for $q = k - 1$

$$\mathcal{V}_k^*(\varphi_l, \mathbf{M}_k) \leq \mathcal{V}_q^*(\varphi_l, \mathbf{M}_k), \quad 0 \leq l \leq T. \quad (75)$$

Finally, by recursively applying (75), the theorem is proved.

REFERENCES

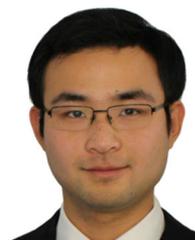
- [1] M. Steinberg, J. Stack, and T. Paluszkiwicz, "Long duration autonomy for maritime systems: Challenges and opportunities," *Auton. Robots*, vol. 40, no. 7, pp. 1119–1122, 2016.
- [2] L. E. Parker, "Path planning and motion coordination in multiple mobile robot teams," in *Proc. Encyclopedia Complexity Sys. Sci.*, New York, NY, USA: Springer, 2009, pp. 5783–5800.
- [3] G. Foderaro, S. Ferrari, and T. A. Wettergren, "Distributed optimal control for multi-agent trajectory optimization," *Automatica*, vol. 50, no. 1, pp. 149–154, 2014.
- [4] S. Ferrari, G. Foderaro, P. Zhu, and T. A. Wettergren, "Distributed optimal control of multiscale dynamical systems: A tutorial," *IEEE Control Syst. Mag.*, vol. 36, no. 2, pp. 102–116, Apr. 2016.
- [5] K. Rudd, G. Foderaro, P. Zhu, and S. Ferrari, "A generalized reduced gradient method for the optimal control of very-large-scale robotic systems," *IEEE Trans. Robot.*, vol. 33, no. 5, pp. 1226–1232, Oct. 2017.
- [6] Y. Yang, R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang, "Mean field multi-agent reinforcement learning," 2018, *arXiv:1802.05438*.

- [7] A. Khan, C. Zhang, D. D. Lee, V. Kumar, and A. Ribeiro, "Scalable centralized deep multi-agent reinforcement learning via policy gradients," 2018, *arXiv:1805.08776*.
- [8] M. Hüttenrauch *et al.*, "Deep reinforcement learning for swarm systems," *J. Mach. Learn. Res.*, vol. 20, no. 54, pp. 1–31, 2019.
- [9] P. Hernandez-Leal, B. Kartal, and M. E. Taylor, "A survey and critique of multiagent deep reinforcement learning," *Auton. Agents Multi-Agent Syst.*, vol. 33, no. 6, pp. 750–797, 2019.
- [10] R. Carmona, M. Laurière, and Z. Tan, "Linear-quadratic mean-field reinforcement learning: Convergence of policy gradient methods," 2019, *arXiv:1910.04295*.
- [11] R. Carmona, M. Laurière, and Z. Tan, "Model-free mean-field reinforcement learning: Mean-field MDP and mean-field Q-learning," 2019, *arXiv:1910.12802*.
- [12] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," 2019, *arXiv:1911.10635*.
- [13] W. Ren, "Represented value function approach for large scale multi agent reinforcement learning," 2020, *arXiv:2001.01096*.
- [14] C. Ferrari, E. Pagello, J. Ota, and T. Arai, "Multirobot motion coordination in space and time," *Robot. Auton. Syst.*, vol. 25, no. 3/4, pp. 219–229, 1998.
- [15] M. Bennewitz, W. Burgard, and S. Thrun, "Finding and optimizing solvable priority schemes for decoupled path planning techniques for teams of mobile robots," *Robot. Auton. Syst.*, vol. 41, no. 2/3, pp. 89–99, 2002.
- [16] M. Yogeswaran and S. Ponnambalam, "Swarm robotics: An extensive research review," in *Proc. Adv. Know. App. Practice*. London, U.K.: IntechOpen, 2010.
- [17] M. Rubenstein, A. Cornejo, and R. Nagpal, "Programmable self-assembly in a thousand-robot swarm," *Science*, vol. 345, no. 6198, pp. 795–799, 2014.
- [18] S. Bandyopadhyay, S.-J. Chung, and F. Y. Hadaegh, "Probabilistic swarm guidance using optimal transport," in *Proc. IEEE Conf. Control Appl.*, 2014, pp. 498–505.
- [19] V. Krishnan and S. Martínez, "Distributed optimal transport for the deployment of swarms," in *Proc. IEEE Conf. Decis. Control*, 2018, pp. 4583–4588.
- [20] L. Bayındır, "A review of swarm robotics tasks," *Neurocomputing*, vol. 172, pp. 292–321, 2016.
- [21] N. Nedjah and L. S. Junior, "Review of methodologies and tasks in swarm robotics towards standardization," *Swarm Evol. Comput.*, vol. 50, 2019, Art no. 100565.
- [22] M. Huang *et al.*, "Large population stochastic dynamic games: Closed-loop Mckean–Vlasov systems and the Nash certainty equivalence principle," *Commun. Inf. Syst.*, vol. 6, no. 3, pp. 221–252, 2006.
- [23] M. Huang, P. E. Caines, and R. P. Malhamé, "The NCE (mean field) principle with locality dependent cost interactions," *IEEE Trans. Autom. Control*, vol. 55, no. 12, pp. 2799–2805, Dec. 2010.
- [24] S. Ferrari and R. F. Stengel, "An adaptive critic global controller," in *Proc. Amer. Control Conf.*, 2002, pp. 2665–2670.
- [25] S. Ferrari and R. Stengel, "Model-based adaptive critic designs," in *Proc. Learn. Approx. Dynam. Program.*. Hoboken, NJ, USA: Wiley-IEEE Press, 2004.
- [26] F. Ramos and L. Ott, "Hilbert maps: Scalable continuous occupancy mapping with stochastic gradient descent," *Int. J. Robot. Res.*, vol. 35, no. 14, pp. 1717–1730, 2016.
- [27] P. Zhu, S. Ferrari, J. Morelli, R. Linares, and B. Doerr, "Scalable gas sensing, mapping, and path planning via decentralized hilbert maps," *Sensors*, vol. 19, no. 7, 2019, Art. no. 1524.
- [28] P. Zhu, H. Wei, W. Lu, and S. Ferrari, "Multi-kernel probability distribution regressions," in *Proc. Int. Joint Conf. Neural Netw.*, 2015.
- [29] P. Zhu, J. Morelli, and S. Ferrari, "Value function approximation for the control of multiscale dynamical systems," in *Proc. IEEE 55th Conf. Decis. Control*, 2016, pp. 5471–5477.
- [30] B. Luo, D. Liu, T. Huang, and D. Wang, "Model-free optimal tracking control via critic-only Q-learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 27, no. 10, pp. 2134–2144, Oct. 2016.
- [31] C. Villani, *Topics in Optimal Transportation*. Providence, RI, USA: Amer. Math. Soc., 2003.
- [32] Y. Chen, T. T. Georgiou, and A. Tannenbaum, "Optimal transport for Gaussian mixture models," *IEEE Access*, vol. 7, pp. 6269–6278, 2018.
- [33] V. M. Panaretos and Y. Zemel, "Statistical aspects of Wasserstein distances," *Annu. Rev. Statist. Appl.*, vol. 6, pp. 405–431, 2019.
- [34] G. Auricchio, F. Bassetti, S. Gualandi, and M. Veneroni, "Computing Kantorovich–Wasserstein distances on d -dimensional histograms using $(d + 1)$ -partite graphs," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 5793–5803.
- [35] M. Kristan, A. Leonardis, and D. Škočaj, "Multivariate online kernel density estimation with Gaussian kernels," *Pattern Recognit.*, vol. 44, no. 10/11, pp. 2630–2642, 2011.
- [36] M. L. Fredman and R. E. Tarjan, "Fibonacci heaps and their uses in improved network optimization algorithms," *J. ACM*, vol. 34, no. 3, pp. 596–615, 1987.
- [37] M. Thorup, "Integer priority queues with decrease key in constant time and the single source shortest paths problem," *J. Comput. Syst. Sci.*, vol. 69, no. 3, pp. 330–353, 2004.



Pingping Zhu (Member, IEEE) received the B.S. degree in electronics and information engineering and the M.S. degree from the Institute for Pattern Recognition and Artificial Intelligence, Huazhong University of Science and Technology, Wuhan, China, in 2006 and 2008, respectively, and the Ph.D. degree in electrical and computer engineering from the University of Florida, Gainesville, FL, USA, in 2013.

He is currently an Assistant Professor with the Department of Computer Science and Electrical Engineering, Marshall University, Huntington, WV, USA. He was a Research Associate with the Department of Mechanical and Aerospace Engineering, Cornell University, Ithaca, NY, USA. His research interests include approximate dynamic programming, reinforcement learning, signal processing, information theoretical learning, machine learning, artificial intelligence, and neural networks.



Chang Liu received the B.S. degrees (double major) in electrical engineering and applied mathematics from Peking University, Beijing, China, in 2011, the M.S. degrees in mechanical engineering and computer science, and the Ph.D. degree in mechanical engineering from the University of California, Berkeley, CA, USA, in 2014, 2016, and 2017, respectively.

He is currently a Postdoctoral Associate with the Sibley School of Mechanical and Aerospace Engineering, Cornell University, Ithaca, NY, USA. His research interests include planning and decision-making of robots, multiagent systems, state estimation and prediction, computer vision, and human–robot collaboration.



Silvia Ferrari (Senior Member, IEEE) received the B.S. degree from Embry–Riddle Aeronautical University, Daytona Beach, FL, USA, and the M.A. and Ph.D. degrees from Princeton University, Princeton, NJ, USA.

She is currently a John Brancaccio Professor of Mechanical and Aerospace Engineering with Cornell University, Ithaca, NY, USA. She was a Professor of Engineering and Computer Science with Duke University, Durham, NC, USA, and the Founder and Director of the NSF Integrative Graduate Education and Research Traineeship and the Fellowship program on Wireless Intelligent Sensor Networks. Her research interests include robust adaptive control of aircraft, learning and approximate dynamic programming, and optimal control of mobile sensor networks.

Dr. Ferrari is a Fellow of American Society of Mechanical Engineers, an Associate Fellow of the American Institute of Aeronautics and Astronautics, and a Member of Society of Photo-Optical Instrumentation Engineers and Society for Industrial and Applied Mathematics. She was the recipient of the Office of Naval Research Young Investigator Award in 2004, the National Science Foundation CAREER Award in 2005, and the Presidential Early Career Award for Scientists and Engineers Award, in 2006.