arXiv:2109.04927v3 [cs.RO] 6 Dec 2021

Learning to Swarm with Knowledge-Based Neural Ordinary Differential Equations

Tom Z. Jiahao*, Lishuo Pan* and M. Ani Hsieh

Abstract-Understanding decentralized dynamics from collective behaviors in swarms is crucial for informing robot controller designs in artificial swarms and multiagent robotic systems. However, the complexity in agent-to-agent interactions and the decentralized nature of most swarms pose a significant challenge to the extraction of single-robot control laws from global behavior. In this work, we consider the important task of learning decentralized single-robot controllers based solely on the state observations of a swarm's trajectory. We present a general framework by adopting knowledge-based neural ordinary differential equations (KNODE) - a hybrid machine learning method capable of combining artificial neural networks with known agent dynamics. Our approach distinguishes itself from most prior works in that we do not require action data for learning. We apply our framework to two different flocking swarms in 2D and 3D respectively, and demonstrate efficient training by leveraging the graphical structure of the swarms' information network. We further show that the learnt singlerobot controllers can not only reproduce flocking behavior in the original swarm but also scale to swarms with more robots.

I. INTRODUCTION

Many natural swarms exhibit mesmerizing collective behaviors, and have fascinated researchers over the past decade [1], [2], [3], [4], [5]. A leading question is how do these global behaviors emerge from local interactions. Such fascination has led to much developments in artificial swarms and multi-agent robotic systems to emulate the swarms in nature. [6], [7], [8]. Central to these developments is the task of single-robot swarm controller synthesis, which has enabled deployment of robot swarms that respects task specifications and real-world constraints.

Some of the earliest works on developing swarm controllers rely heavily on physical intuitions and design controllers in a bottom-up fashion. Boids was developed by combining rules of cohesion, alignment, and separation to mimic the flocking behavior in natural swarms [6]. Selfdriven particles were used to model the emergence of collective behaviors in biologically motivated swarms [9]. Flocking controllers with provably correct stability guarantees have also been developed for swarms with fixed and dynamic communication network topologies [7], [10]. These early works laid the foundation of decentralized swarm control and offered a glimpse of the myriad of possible swarm behaviors achievable using local single-agent controllers. In recent years, deep learning has enabled pattern discovery from complex and high-dimensional data sets. The use of neural networks (NNs) have shown promising results in a wide range of applications owing to their expressive power. This has opened up potential avenues for data-driven learning of single-robot swarming control strategies in more efficient and scalable ways. In this work, we leverage recent advances in scientific machine learning and employ knowledge-based neural ordinary differential equations (KNODE) [11] for learning swarm controllers directly from observations of a swarm. We demonstrate that through our top-down approach to controller synthesis, global behaviors of different swarms can be successfully reproduced based on the past observations of their evolution.

II. RELATED WORKS

Various data-driven methods have been used to model local control policy in swarms. Feedforward neural networks have been used to approximate decentralized control policies by training on the observation-action data from a global planner [12]. Furthermore, deep neural networks have been used to model higher order residual dynamics to achieve stable control in a swarm of quadrotors [13]. Recently, graph neural networks (GNN) have been extensively used in swarms, owing to their naturally distributed architecture. GNN allows efficient information propagration through networks with underlying graphical structures [14], and have been noted for their stability and permutation equivariance [15]. Decentralized GNN controllers have been trained with global control policies to imitate swarm behaviors [14], [16]. All these works pose the controller synthesis problem as an imitation learning problem, and require knowledge of the actions resulting from an optimal control policy for learning or improving the local controllers. In practice, action data can be difficult to access, especially when learning behaviors from natural or adversarial swarms. In addition, GNNs can potentially allow a robot to access the state information of robots outside its communication range through information propagation. The true extent of decentralization may therefore be limited when more propagation hops are allowed.

Deep reinforcement learning has also been applied to swarms for various applications [17]. Early works like [18] learn a decentralized control policy for maintaining distances within a swarm and target tracking. An inverse reinforcement learning algorithm was presented in [19] to train a decentralized policy by updating the reward function alongside the control policy based on an expert behavior. In addition, GNNs have also been used within the reinforcement learning

This work was supported by ARL DCIST CRA W911NF-17-2-0181 and Office of Naval Research (ONR) Award No. 14-19-1-2253.

All authors in this work are with the GRASP Laboratory, University of Pennsylvania, Philadelphia, PA 19104, USA. {zjh, panls, m.hsieh}@seas.upenn.edu

^{*}Equal contribution.

framework for learning connectivity for data distribution [20]. However, reinforcement learning is usually employed to solve task-specific problems with well-defined goals and need to tackle the challenge of reward shaping. The specific objectives of swarms may be difficult to discern from only observations, and therefore reinforcement learning is often not suitable for learning global behaviors from solely observational data.

The contribution of this work is three-fold. First, we demonstrate the feasibility of learning single-robot controllers that can achieve the observed global swarming behaviors from only swarm trajectory data. Second, we propose a generalized model for incorporating known robot dynamics to facilitate learning single-robot controllers. Lastly, we show how to efficiently scale KNODE for learning from local information in a multi-agent setting.

III. PROBLEM FORMULATION

We consider the problem of learning single-robot controllers based on the observations of the trajectory of a swarm. We assume that the swarm is homogeneous, *i.e.*, all robots in the swarm use the same controller. Given a swarm of n agents, we make m observations at sampling times $T = \{t_1, t_2, ..., t_m\}, t_i \in \mathbb{R}$ given by

$$\begin{bmatrix} \mathbf{Z}^{T}(t_{1}) \\ \mathbf{Z}^{T}(t_{2}) \\ \vdots \\ \mathbf{Z}^{T}(t_{m}) \end{bmatrix} = \begin{bmatrix} \mathbf{z}_{1}(t_{1}) & \mathbf{z}_{2}(t_{1}) & \cdots & \mathbf{z}_{n}(t_{1}) \\ \mathbf{z}_{1}(t_{2}) & \mathbf{z}_{2}(t_{2}) & \cdots & \mathbf{z}_{n}(t_{2}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{z}_{1}(t_{m}) & \mathbf{z}_{2}(t_{m}) & \cdots & \mathbf{z}_{n}(t_{m}) \end{bmatrix},$$

where the matrix $\mathbf{Z}(t_i) \in \mathbb{R}^{n \times d}$ is the observations of the states of all *n* agents at t_i , and the vector $\mathbf{z}_i(t_j) \in \mathbb{R}^d$ is the state of agent *i* observed at t_j with dimension *d*. For instance, in a first-order system, an agent modeled as a rigid body in a 3-dimensional space has d = 6, where the first three dimensions correspond to the positions and the last three the orientations. Our goal is to learn a single-robot controller solely from the observations \mathbf{Z} . Notice that control inputs are not assumed to be part of the observations.

The evolution of each individual robot's state can be described by the true dynamics given by

$$\dot{\mathbf{z}}_i(t) = f_i(\mathbf{z}_i, u_i),\tag{1}$$

where \mathbf{z}_i is the state of robot *i*, and u_i is its control law. The function $f_i(\cdot, \cdot)$ defines the dynamics given the state of robot and control law u_i . It is assumed that all robots in the swarm have the same dynamics and control strategy, and therefore we can drop the subscripts and rewrite (1) as $\dot{\mathbf{z}}_i(t) = f(\mathbf{z}_i, u)$ for all *i*. The control law *u* is a function of the states of other robots in the swarm, and defines the interaction between robot *i*. For example, a communication radius may be enforced by the control law *u* to let each robot only interact with its neighbors.

The dynamics of the entire swarm can be written as a collection of the single-robot dynamics as

$$\dot{\mathbf{Z}}(t) = [\dot{\mathbf{z}}_1(t), \dot{\mathbf{z}}_2(t), \cdots, \dot{\mathbf{z}}_n(t)]^T.$$
(2)

Given the initial conditions of all robots \mathbf{Z}_0 at t_0 , the states of all robots at t_1 is given by

$$\mathbf{Z}(t_1) = \mathbf{Z}_0 + \int_{t_0}^{t_1} \dot{\mathbf{Z}}(t) dt.$$
(3)

In practice, the integration in (3) is performed numerically. Our task is to find a single-robot control law parameterized by θ as part of the single-robot dynamics given by

$$\dot{\mathbf{z}}_i(t) = \hat{f}(\mathbf{z}_i, \hat{u}_{\boldsymbol{\theta}}), \tag{4}$$

where \hat{u}_{θ} is the single-robot control law parameterized by θ . The learnt controller should best reproduce the observed global swarm behaviors. Note that the high dimensionality of a swarming system means that similar collective dynamics can be achieved with very disparate collections of single-robot trajectories. This suggests that it may be impractical to predict each individual trajectory in a swarm over long time horizons. Instead, we focus on learning and reproducing the global behaviors of swarms based on metrics, which we will formalize in later sections.

IV. KNOWLEDGE-BASED NEURAL ORDINARY DIFFERENTIAL EQUATIONS (KNODE)

KNODE is a scientific machine learning framework that applies to a general class of dynamical systems. It has been shown to model a wide variety of systems with nonlinear and chaotic dynamics, with robustness to noise and irregularly sampled data [11]. In our problem, we assume a singlerobot dynamics in the form of (4). From a dynamical systems perspective, $\hat{f}(\mathbf{z}_i, \hat{u}_{\theta})$ is a vector field. This makes KNODE a suitable method to learn $\hat{f}(\mathbf{z}_i, \hat{u}_{\theta})$ because it directly models vector fields using neural networks [11]. To put KNODE in the context of our learning problem, given some known swarm dynamics $\tilde{f}(\mathbf{Z})$ as knowledge, KNODE optimizes for the control law as part of a dynamics given by

$$\dot{\mathbf{z}}_i(t) = \hat{f}(\mathbf{z}_i, \hat{u}_{\boldsymbol{\theta}}, \tilde{f}(\mathbf{Z})), \tag{5}$$

where the control law \hat{u}_{θ} is a neural network, and \hat{f} defines the coupling between the knowledge and the rest of the dynamics. While the original KNODE linearly couples a neural network with \tilde{f} using a trainable matrix [11], we note that the way knowledge gets incorporated is flexible. In later sections we will demonstrate how to effectively incorporate knowledge for learning single-robot controllers. Furthermore, the ability to incorporate knowledge will require less training data [21], [11].

We minimize the mean squared error (MSE) between the observed trajectories and the trajectories predicted from the estimate dynamics using \hat{u}_{θ} for robot *i*. A loss function is given by

$$L(\boldsymbol{\theta}) = \frac{1}{m-1} \sum_{j=1}^{m-1} \sum_{i=1}^{n} \|\hat{\mathbf{z}}_i(t_{j+1}, \mathbf{z}_i(t_j)) - \mathbf{z}_i(t_{j+1})\|_2^2,$$
(6)

where $\hat{\mathbf{z}}_i(t_{j+1}, \mathbf{z}_i(t_j))$ is the estimated state of robot *i* at t_{j+1} generated using the initial condition $\mathbf{z}_i(t_j)$ at t_j , and

it's given by

$$\hat{\mathbf{z}}_i(t_{j+1}, \mathbf{z}_i(t_j)) = \mathbf{z}_i(t_j) + \int_{t_j}^{t_{j+1}} \hat{f}(\mathbf{z}_i, \hat{u}_{\boldsymbol{\theta}}, \tilde{f}(\mathbf{Z})) dt.$$
(7)

Intuitively, the loss function in (6) computes the one-stepahead estimated state of all robots from every snapshot in the observed trajectory, and then computes the average MSE between the estimated and observed states for the entire trajectory from t_1 to t_{m-1} .

Our learning task can then be formulated as an optimization problem given by

$$\min_{\boldsymbol{\theta}} \quad L(\boldsymbol{\theta}), \tag{8}$$

s.t.
$$\dot{\mathbf{z}}_i = \hat{f}(\mathbf{z}_i, \hat{u}_{\theta}, \hat{f}(\mathbf{Z})), \text{ for all } i,$$
 (9)

which includes the dynamics constraint for all robots in the swarm. The parameters θ can then be estimated by $\theta = \arg \min_{\theta} L(\theta)$. The gradients of θ with respect to the loss can be computed by either the conventional backpropagation or the adjoint sensitivity method. The adjoint sensitivity method has been noted as a more memory efficient approach than backpropagation, though at the cost of training speed [22]. In this work, we use the adjoint method for training similar to that in [23] and [11].

V. METHOD

In this section, we walk through the process for constructing $\hat{f}(\mathbf{z}_i, \hat{u}_{\theta}, \tilde{f}(\mathbf{Z}))$ in the context of learning to swarm and the incorporation of knowledge in the form of known single robot dynamics.

A. Decentralized Information Network

We assume a robot in a swarm can only use its local information as inputs to its controller. To incorporate this assumption, we impose a decentralized information network on the swarm. Specifically, we assume robots have finite communication radii as denoted by d_{cr} . In addition, each robot can only communicate with a maximum number of neighbors, including itself, as denoted by k. We refer to the robots within this radius as the *active neighbors*. If there are more than k neighbors within a robot's communication radius, the closest k neighbors are considered to be *active*.

We leverage the communication graph of the swarm to compute the local information for each robot at each time step. The communication graph at time t can be described by a graph shift operator $\mathbf{S}(t) \in \mathbb{R}^{n \times n}$, which is a binary adjacency matrix computed based on d_{cr} and the positions of all robots at each time step. In this work, we treat the communication radius d_{cr} as a hyperparameter. Note that the communication graph is time-varying because the information network changes as robots move around in a swarm. Then $\mathbf{S}_{ij}(t) = 1$ if the Euclidean distance between agents i and j is less than or equal to d_{cr} , and $\mathbf{S}_{ij}(t) = 0$ otherwise. The index set of the neighbors of robot i at time t is therefore given by

$$\mathcal{N}_i(t) = \{ j | j \in \mathcal{I}, \mathbf{S}_{ij}(t) = 1 \}, \tag{10}$$

where $\mathcal{I} = 1, ..., n$ is the index set of all robots. Note that set of neighbors of robot *i* also includes itself. At time *t*, the information kept by robot *i* is the matrix $\mathbf{Y}_i(t) \in \mathbb{R}^{k \times d}$ given by

$$\mathbf{Y}_{i}(t) = g(\{\mathbf{z}_{j}(t) | j \in \mathcal{N}_{i}(t)\}, k),$$
(11)

where the function $g(\cdot, k)$ maintains the dimension of the matrix $\mathbf{Y}_i(t)$, and forms the rows of matrix $\mathbf{Y}_i(t)$ using the state information of robot *i*'s active neighbors in ascending order of their Euclidean distance from robot *i*. Naturally, robot *i*'s state is always in the first row because its distance to itself is 0. If there are fewer than *k* active neighbors within a robot's communication radius, the remaining rows in $\mathbf{Y}_i(t)$ are padded with zeros. In this work, *k* is treated as a hyperparameter.

The matrix $\mathbf{Y}(t)$ represents the local information accessible to each robot at time t and it completes the decentralized information network of the swarm. In summary, (10) and (11) enforces the assumptions of finite communication and perception radii for each robot.

B. Information Time Delay

In addition to a decentralized information structure, we further assume that each robot only gets delayed state information from its neighboring robots by a time lag τ . This is to emulate the latency in agent communication in real swarms. With time delay, the information accessible to robot *i* in (11) becomes

$$\mathbf{Y}_{i}(t) = \begin{bmatrix} \mathbf{z}_{i}^{T}(t) \\ g\left(\{\mathbf{z}_{j}(t-\tau) | i \neq j, j \in \mathcal{N}_{i}(t-\tau)\}, k-1\right) \end{bmatrix}.$$
(12)

Fig. 1 shows an example of the information structure described by (12) using k = 3. The process of constructing $\mathbf{Y}_i(t)$ for all $t \in T$ in (10), (11) and (12) leverages the graphical structure of the swarm's information network. During training, the collection of delayed neighbor information is done efficiently through the matrix multiplication $\mathbf{S}(t-\tau)\mathbf{Z}(t-\tau)$, which leaves for each robot only the state information of its neighbors at $t - \tau$. Then for robot *i* we append the *i*th row of $[\mathbf{S}(t-\tau)\mathbf{Z}(t-\tau)]$ to its own state $\mathbf{z}_i^T(t)$. Finally we only keep k rows of the resulting matrix to form $\mathbf{Y}_{i}(t)$. Compared to some GNN approaches [14], [15], the information structure $\mathbf{Y}_{i}(t)$ in our work is more explicit. A robot with GNN controllers can only access the diffused state information from other robots, *i.e.* the neighbors' information has been repeatedly multiplied by the graph operators before reaching this robot. In this work, we directly let each robot access the state information of its active neighbors. In real-world implementation of robot swarms, our proposed information structure in (12) is more realistic as each robot can easily subscribe to or observe its neighbors' states. In addition, the information structure $\mathbf{Y}_{i}(t)$ enables scalable learning as we can treat the robots in a swarm as batches. As a result, training memory scales linearly with the number of robots in the swarm, and training speed scales sub-linearly.



Fig. 1. Decentralized information network for robot 0 with time delay τ , and 3 active neighbors. The image shows robot 0's egocentric view, where 8 neighbors are within its communication range d_{cr} . Only the closest three neighbors contribute to the information structure of robot 0. Their states from $t - \tau$ are ordered based on their proximity to robot 0 to form $\mathbf{Y}_0(t)$.

C. Knowledge Embedding

In this work, a potential-function-based obstacle avoidance strategy similar to [24] is used as knowledge. Let the distance between robot *i* and an obstacle \mathcal{O} be $d_{\mathcal{O}}(\mathbf{z}_i)$, where \mathbf{z}_i is the state and includes the position of robot *i*. The potential function is then given by

$$U_{\mathcal{O}}(\mathbf{z}_i) = \begin{cases} \frac{\lambda}{2} \frac{1}{d_{\mathcal{O}}^2(\mathbf{z}_i)} & \text{if } d_{\mathcal{O}}(\mathbf{z}_i) \le d_0, \\ 0 & \text{otherwise,} \end{cases}$$
(13)

where λ is the gain, and d_0 is the obstacle influence threshold (*i.e.* the distance within which the potential function becomes active). Based on this potential function, the repulsive force to avoid the obstacle \mathcal{O} is given by

$$F_{\mathcal{O}}(\mathbf{z}_i) = \begin{cases} -\nabla U_{\mathcal{O}}(\mathbf{z}_i) & \text{if } d_{\mathcal{O}}(\mathbf{z}_i) \le d_0, \\ 0 & \text{otherwise}, \end{cases}$$
(14)

When multiple obstacles are present, the repulsive forces computed from each obstacle are summed for a resultant repulsive force. For collision avoidance, we assume that each agent will only actively avoid its closest neighbor within d_0 at any given time.

Assuming that the robots in a swarm follow first-order dynamics, we combine the decentralized information network in (11) and the knowledge in (14) into the dynamics given by

$$\dot{\mathbf{z}}_{i} = \hat{f}(\mathbf{z}_{i}, \hat{u}_{\boldsymbol{\theta}}(\mathbf{Y}_{i}, \mathbf{z}_{i})) - \sum_{j} \lambda_{j} \nabla U_{\mathcal{O}_{j}}(\mathbf{z}_{i}), \quad (15)$$

where u_{θ} is a neural network, and λ_j is a trainable gain for avoiding obstacle \mathcal{O}_j . Note that Eqn. (15) further illustrates how our framework differs from imitation learning. While \dot{z}_i and \hat{u}_{θ} are the learnt dynamics and control policy which drive the system, they do not have to be part of the training data.

VI. LEARNING TO FLOCK IN 2D

We first use a global controller proposed by [7] to generate observations for our learning problem.

A. Simulation in 2D and training

This global controller achieves stable flocking, which ensures eventual velocity alignment, collision avoidance and group cohesion in a swarm of robots. The robots follow the double integrator dynamics given by

$$\mathbf{r}_i = \mathbf{v}_i, \dot{\mathbf{v}}_i = \mathbf{u}_i, \quad i = 1, ..., n,$$
(16)

where \mathbf{r}_i is the 2D position vector of robot *i*, \mathbf{v}_i is its velocity vector. The full state of each robot is therefore $\mathbf{x} = [\mathbf{r}, \mathbf{v}] \in \mathbb{R}^4$. The control law \mathbf{u}_i is given by

$$\mathbf{u}_{i} = -\sum_{j \in \mathcal{N}_{i}} (\mathbf{v}_{i} - \mathbf{v}_{j}) - \sum_{j \in \mathcal{N}_{i}} \nabla_{\mathbf{r}_{i}} V_{ij}, \qquad (17)$$

where V_{ij} is a differentiable, nonnegative, and radially unbounded function of the distance between robot *i* and *j* [7]. The first summation term in (17) aims to align the velocity vector of robot *i* with those of its flockmates, while the second summation term is the total potential field around robot *i* responsible for both collision avoidance and cohesion [7]. The set N_i is the set of all robots in the swarm for the global controller.

We use the explicit fourth-order Runge-Kutta method to simulate the dynamics in (17) with a step time of 0.01. Given n robots, their locations are initialized uniformly on a disk with radius \sqrt{n} to normalize the density within the swarm. The velocities of robots are initialized uniformly with magnitudes between [0, 3]. Additionally, a uniformly sampled velocity bias with magnitude between [0, 3] is added to the swarm. A total of 50 trajectories are simulated, each with a total of 2000 steps. The lengths of the trajectories is chosen such that the swarms will converge to stable flocking. We use 30 trajectories as the training data, and the remaining 20 as the testing data. We added zero-mean Gaussian noise with variance 0.001 to the training trajectories. This is known as *stabilization noise* in modeling dynamical systems and has been shown to improve model convergence [25].

The training model follows (15). There are no obstacles to avoid in the 2D case, so the potential function is only used to avoid collision among the agents. Specifically, we let each robot avoid its closest neighbor at every time step. For the controller \hat{u}_{θ} , we use a one layer neural network with 128 hidden units, and a hyperbolic tangent activation function. The trainable gain for collision avoidance is defined as $\lambda = a + \phi^2$, where a is a positive number for setting the minimum amount of force to avoid collision. The single parameter ϕ is trained together with the neural network. We do not assume information delay in the 2D case.

B. Evaluating flocking in 2D

We evaluate 2D flocking behavior using two metrics: **Average velocity difference** (*avd*) measures how well the velocities of robots are aligned. It is given by

$$avd(t) = \frac{2}{n(n-1)} \sum_{i \neq j} ||\mathbf{v}_i(t) - \mathbf{v}_j(t)||_2.$$
 (18)

Average minimum distance to a neighbor (*amd*) measures the cohesion between agents in both 2D and 3D when flocking is achieved. It is given by

$$amd(t) = \frac{1}{n} \sum_{i=1}^{n} \min_{j} ||\mathbf{r}_i(t) - \mathbf{r}_j(t)||_2.$$
 (19)

amd should decrease as the robots move closer together, but it should not reach zero if collision avoidance is in place. To generate trajectories using the learnt controller, we use it to replace (17) in the dynamics described by Eqn. (16) for acceleration control.

C. 2D Results

Fig. 2 shows four snapshots of the swarm trajectory generated using the trained single-robot controller, and provides a qualitative comparison between the prediction and ground truth. The controller used to produce these snapshots were trained with $d_{cr} = 5$ and k = 6. The robots are initialized using the initial states from the testing trajectory. It can be observed that the predicted swarm achieves velocity alignment while the robots stay apart from each other, indicating the emergence of flocking behavior. This can be further verified by the metrics for 2D flocking as shown in Fig. 3. The predicted swarm trajectory follows similar trends as the ground truth under both metrics.

Furthermore, we deployed the trained controller on larger swarms to test its scalability. Each of these swarms are uniformly initialized in a ball around the origin, with the same robot density as the training data. Fig. 5 shows the controller performance on swarms of sizes from 10 to 90. It can be observed that amd remains largely consistent, demonstrating that collision avoidance is effective and cohesion is in place even as the swarm size increases. Although avd degrades as the swarm size increases, it remains low enough that some velocity alignment is achieved. As a qualitative illustration, Fig. 4 shows six snapshots of a swarm of 100 robots using the learnt controller. Although qualitatively velocity alignment can be observed in the predicted trajectories from the snapshots, the global behavior is different from the simulation. This is because the simulation uses the global controller while our prediction uses the decentralized controller learnt from the 10-agent data. In other words, the predictions are the best effort to mimic the centralized 100-agent swarm using the learnt decentralized controller. We do note that with some initialization, the predicted 100-agent swarm tends to split into subswarms. This is not unexpected since stability of the original controller is only guaranteed under certain conditions [7], [10].

We further conducted analysis on the hyperparameters d_{cr} and k with respect to 2D flocking. Grid searches are performed on both *avd* and *amd* by varying d_{cr} and k. For each grid, the average of the last 10 steps of a 2000-step trajectory are computed for 20 different initial conditions. The average over these 20 different initial conditions is then reported in the grid. It can be observed from Fig. 6 that the *avd* is poor for both small values of d_{cr} and k, while *amd*



Fig. 2. Predicted trajectory of 10 robots using the learnt controller ($d_{cr} = 5, k = 6$) with the same initial states as the testing trajectory (ground truth). The subfigures (a)(b)(c)(d) show the snapshots of the swarm at t = 0, 100, 600, and 1200 respectively.



Fig. 3. The metrics for the learnt 2D controller ($d_{cr} = 5, k = 6$) show (a) average velocity difference, and (b) average minimum distance to a neighbor. The 95% confidence intervals are based on 20 sets of testing trajectories.

is largely affected by d_{cr} only. This grid search result agrees with intuition and can help with hyperparameter selection.

VII. LEARNING TO FLOCK IN 3D

Next, we apply our learning method on the 3D simulation of boids. Boids was introduced to emulate flocking behaviors and led to the creation of *artificial life* in the field of computer graphics [6]. The flocking behavior of boids is more challenging to learn because (1) they have higher dimensionality, and (2) their steady state flocking behavior is more complex than the 2D flocking in the previous section when the swarm is confined within limited volume.

A. Simulation in 3D and training

Boids are simulated based on three rules:

- **cohesion** each boid moves towards the average position of its neighboring boids.
- **alignment** each boid steer towards the average heading of its neighboring boids.
- **separation** each boid steer towards direction with no obstacles to avoid colliding into its neighboring boids.

While cohesion and collision avoidance are grouped into one term in the 2D flocking case, boids use two separate terms. Furthermore, the boids in simulation are confined in a cubic space and are tasked to avoid the boundaries.



Fig. 4. Predicted trajectory of 100 robots using the learnt controller ($d_{cr} = 5, k = 6$) with uniformly initialized positions. The subfigures (a)(b)(c)(d)(e)(f) show the snapshots of the swarm at t = 0,200,400,800,1000 and 1200 respectively.



Fig. 5. Box plot of (a) average velocity difference (avd), and (b) average minimum distance to a neighbor (amd) on scaling to different swarm sizes using a trained controller in 2D. For each swarm size, the box represents the statistics of 15 runs using different initial conditions.

Boids are simulated in Unity [26]. We follow the default settings with a minimum boids speed of 2.0, a maximum speed of 5.0, a communication radius of 2.5 (ball), a collision avoidance range of 1.0, a maximum steering force of 3.0, and the weights of cohesion, alignment, and separation steering force are all set to 1.0. For obstacle avoidance we set the scout sphere radius as 0.27, the maximum search distance as 5.0, and the weight of obstacle avoidance steering force as 10.0. Boids are simulated in a cubic space with an equal side length of 10, with each axis ranging from -5 to 5. The boids' positions are randomly initialized within a sphere of



Fig. 6. Grid search on (a) the average velocity difference (avd) and (b) average minimum distance (amd) to a neighbor using different communication radii and number of active neighbors. The grid values are computed as the average over trajectories using 20 different initial conditions.

radius 5 centered at origin, and their velocities vectors are randomly initialized with a constant magnitude.

Unity can log both the positions and velocities of boids. However, to make the learning task more challenging, we only use the positions and orientations of the boids for training. For a swarm of 10 boids, we simulate 22 trajectories, each with a total of 1700 steps. We discard the first 10 time steps to remove simulation artifacts (There are 'jumps' in the first few steps of simulation) and only use the remaining 1690 steps. We use 2 trajectories for training and the remaining 20 as the testing data. Zero-mean Gaussian noise with variance 0.01 is added to the training trajectories.

The training model follows (15). The controller \hat{u}_{θ} uses a one layer neural network with 128 hidden units and a hyperbolic tangent activation function. In addition to collision avoidance, we also include the knowledge for avoiding the boundaries of the cubic space. This is implemented by treating the closest point on each boundary as an obstacle at any given time. Collision and obstacle avoidance use different gains, both of which are defined as $\lambda = \phi^2$, where ϕ is trained. We further assume an information delay of 1.

B. Evaluating flocking in 3D

Average minimum distance to a neighbor (amd) from (19) is also used for 3D flocking to measure the cohesion between robots. However, avd is not a good metric for evaluating flocking in 3D for two reasons: (1) boids only achieve velocity alignment with the local flockmates because of the presence of obstacles, and (2) boids form subswarms. As a result, global velocity alignment is often not achievable at steady state flocking. We instead compare the Proper orthogonal decomposition (POD) modes of the true and predicted trajectory to check how similar the energy distributions are in their respective dynamics. Built on singular value decomposition, POD is a model order reduction technique for nonlinear high-dimensional dynamical systems. It first decomposes the trajectory of a system into orthonormal modes, and then truncates the system by selecting from these modes to form a low-rank basis that captures the most *energy* of the system [27]. Systems with similar dynamics should have similar distributions of POD modes when their energies are arranged in descending order. To measure the shift in the distribution of POD modes between the predicted trajectories and ground truth, we further employ the Kullback-Leibler divergence (KLD), which measures the statistical distance between probability distributions [28]. Together, we first perform POD on trajectories to find the distribution of their energies. Then we apply KLD on the normalized POD distribution to quantitatively measure the shift in this distribution from the ground truth. We name this metric POD-KLD. To generate trajectory predictions, we directly use the learnt controller for velocity control of the swarm.

C. 3D Results

Fig. 7 shows a qualitative comparison between the testing data and the trajectory generated by a controller trained with $d_{cr} = 2$ and k = 6 using the same initial conditions.



Fig. 7. Predicted trajectories of 10 robots using the learnt controller ($d_{cr} = 2, k = 6$) and the same initialization as the testing trajectory (ground truth). The subfigures (a.i)(a.ii)(a.iii) show snapshots of the ground truth trajectory at t = 0, 400, 800, and (b.i)(b.ii)(b.iii) show the eventual flocking and the formation process of subswarms at t = 0, 400, 800. The light blue lines connect the neighbors in the swarm.



Fig. 8. The metrics for the learnt 3D controller $(d_{cr} = 2, k = 6)$. (a) Average minimum distance to a neighbor of the predicted trajectory converges, and (b) the distribution of the first 10 POD modes of the predictions and ground truth are similar. The 95% confidence intervals of *amd* are based on 20 sets of testing trajectories.

The predicted trajectory shows the formation of subswarms during steady state flocking similar to that of the testing trajectory. Empirically the robots are more likely to form a single swarm at steady state when the robots are initialized closer to each other. The metrics for the learnt controllers are shown in Fig. 8. It can be seen that group cohesion is achieved as both the predicted and true swarm show similar trends for *amd*. Furthermore, the distributions of POD modes between the prediction and testing data are similar, indicating similar dynamics.

We also test the scaling ability of the learnt controller on larger swarms of sizes ranging from 10 to 90. Each of these swarms are uniformly initialized in a ball around the origin, with the same robot density as the training data. Fig. 9 shows the metrics on trajectories of different swarm sizes using the same learnt controller. It can be observed that the trend for *amd* is better than the 2D case as swarm size increases. This can be explained by the fact that robots are confined in a cubic space and do not travel too far from each other. Fig. 10 shows comparison between predictions and simulation when there are 50 robots. Notice that our prediction forms subswarms with this size. This may also occur in simulations of 50 agents in Unity. Both qualitatively



Fig. 9. Box plot of (a) average minimum distance to a neighbor (*amd*), and (b) POD-KLD of trajectories generated by a learnt controller on swarms of different sizes in 3D. For each swarm size, the box represents the statistics of 15 runs using different initial conditions.



Fig. 10. The flocking of 50 robots using the learnt controller $(d_{cr} = 2, k = 6)$ with uniformly initialized positions. The subfigures (a)(b)(c)(d)(e)(f) show the snapshots of the swarm at t = 0,200,400,1300,1600,1900 respectively. The light blue lines connect the neighbors in the swarm.

and quantitatively, the controller learnt in 3D scales better than that in 2D. One possible reason is that the 3D simulation itself is decentralized, while the ground truth controller in 2D is centralized. Hence, the predicted trajectories of a larger swarm in 3D is closer to that in simulation.

A grid search is also performed on the hyperparameters d_{cr} and k for 3D flocking. The results are shown in Fig. 11. While a small k leads to poorer metrics, communication range d_{cr} does not affect the metrics as significantly as in the 2D case. This may be due to fact that the swarm in 3D are confined in a fixed volume, and therefore the higher



Fig. 11. Grid search on (a) the mean of average minimum distance to a neighbor (amd), (b) the median of amd, and (c) the mean of KL divergence of the POD modes using different communication radii and number of active neighbors. The grid values are computed over trajectories using 20 different initial conditions. For amd, the grids in white represent values greater than 3. It can be observed that k, the number of neighbors to keep has large influence on the metrics.

density of robots leads to higher chance for the robots to come within each other's communication range even if their communication range is small. Additionally, it can be seen that not all trained models converge. Especially for small k, cohesion may not be achieved in the resulting swarm. Visual inspections reveal that these instances correspond to when robots overcome the obstacle avoidance potential function and leave the cubic space. Since such singular cases increase the average amd dramatically, to better assess the performance we also plot the medians of amd in Fig. 11. Another observation is that the performance degrades slightly for large k. This can be explained by the increase in the number of neural network parameters - an increment of 1 in k correspond to an increase of 256 parameters as the input size increases. Since the training data and training time are unchanged, a larger neural network may tend to underfit.

VIII. DISCUSSION

Our experiments show that the model proposed in (15) is able to learn flocking in both 2D and 3D using appropriate hyperparameters d_{cr} and k. The choices of d_{cr} and k and the corresponding learnt controllers can inform how the extent of decentralization can affect flocking behavior in robot swarms. Furthermore, we note that the collision avoidance strategy which we used as knowledge does not guarantee collision-free trajectories. This is evident in Fig. 11 where robots using some trained controllers leave the confined box. However, the use of this collision avoidance strategy demonstrates the flexibility of our proposed framework for embedding known knowledge about single-robot dynamics, and users are free to incorporate any knowledge including but not limited to collision avoidance strategies.

IX. CONCLUSION AND FUTURE WORK

We have introduced an effective machine learning algorithm for learning to swarm. Specifically, we applied the algorithm to flocking swarms in 2D and 3D respectively. In both cases, the learnt controllers are able to reproduce global flocking behavior similar to the ground truth. Furthermore, the learnt controllers can scale to larger swarms to produce flocking behaviors. We have shown the effectiveness of knowledge embedding in learning decentralized controllers, and demonstrated the feasibility of learning swarm behaviors from state observations alone, distinguishing our work from prior works on imitation learning. For future work, we plan to learn from real-world data, and implement the learnt controllers on physical robot platforms. In addition, we hope to employ neural networks with special properties to derive stability guarantees for the learnt controllers.

REFERENCES

- A. Okubo, "Dynamical aspects of animal grouping: Swarms, schools, flocks, and herds," Advances in Biophysics, vol. 22, pp. 1–94, 1986.
- [2] G. Flierl, D. Grünbaum, S. Levins, and D. Olson, "From individuals to aggregations: the interplay between behavior and physics," *J. Theor. Biol.*, vol. 196, no. 4, pp. 397–454, Feb. 1999.
- [3] K. Warburton and J. Lazarus, "Tendency-distance models of social cohesion in animal groups," J. Theor. Biol., vol. 150, no. 4, pp. 473– 488, Jun. 1991.

- [4] C. M. Breder, "Equations descriptive of fish schools and other animal aggregations," *Ecology*, vol. 35, no. 3, pp. 361–370, 1954. [Online]. Available: http://www.jstor.org/stable/1930099
- [5] T. Vicsek, "A question of scale," Nature, vol. 411, pp. 421-421, 2001.
- [6] C. W. Reynolds, "Flocks, herds and schools: A distributed behavioral model," in *Proc. of the 14th annual conf. on Comput. graphics and interactive tech.*, 1987, pp. 25–34.
- [7] H. G. Tanner, A. Jadbabaie, and G. J. Pappas, "Stable flocking of mobile agents, part i: Fixed topology," in *42nd IEEE Int. Conf. on Decision and Control*, vol. 2. IEEE, 2003, pp. 2010–2015.
- [8] M. Rubenstein, A. Cornejo, and R. Nagpal, "Robotics. programmable self-assembly in a thousand-robot swarm," *Science*, vol. 345, pp. 795– 9, 08 2014.
- [9] T. Vicsek, A. Czirók, E. Ben-Jacob, I. Cohen, and O. Shochet, "Novel type of phase transition in a system of self-driven particles," *Phys. Rev. Lett.*, vol. 75, pp. 1226–1229, Aug 1995.
- [10] H. Tanner, A. Jadbabaie, and G. Pappas, "Stable flocking of mobile agents, part ii: Dynamic topology," *Departmental Papers (ESE)*, vol. 2, 05 2003.
- [11] T. Z. Jiahao, M. A. Hsieh, and E. Forgoston, "Knowledge-based learning of nonlinear dynamics and chaos," *Chaos*, vol. 31, no. 11, p. 111101, 2021.
- [12] B. Riviere, W. Honig, Y. Yue, and S.-J. Chung, "Glas: Global-to-local safe autonomy synthesis for multi-robot motion planning with endto-end learning," *IEEE Robot. and Automat. Lett.*, vol. 5, no. 3, p. 4249–4256, Jul 2020.
- [13] G. Shi, W. Hönig, Y. Yue, and S.-J. Chung, "Neural-swarm: Decentralized close-proximity multirotor control using learned interactions," in 2020 IEEE Int. Conf. on Robot. and Automat. IEEE, 2020, pp. 3241–3247.
- [14] E. Tolstaya, F. Gama, J. Paulos, G. Pappas, V. Kumar, and A. Ribeiro, "Learning decentralized controllers for robot swarms with graph neural networks," in *Conf. on Robot Learn*. PMLR, 2020, pp. 671–682.
- [15] F. Gama, E. Tolstaya, and A. Ribeiro, "Graph neural networks for decentralized controllers," in *ICASSP 2021-2021 IEEE Int. Conf. on Acoust., Speech and Signal Process.* IEEE, 2021, pp. 5260–5264.
- [16] S. Zhou, M. J. Phielipp, J. A. Sefair, S. I. Walker, and H. B. Amor, "Clone swarms: Learning to predict and control multi-robot systems by imitation," 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Nov 2019.
- [17] M. Hüttenrauch, A. Šošić, and G. Neumann, "Deep reinforcement learning for swarm systems," *Journal of Machine Learning Research*, vol. 20, no. 54, pp. 1–31, 2019.
- [18] M. Hüttenrauch, A. Sosic, and G. Neumann, "Guided deep reinforcement learning for swarm systems," *CoRR*, vol. abs/1709.06011, 2017.
- [19] A. Šošić, W. R. KhudaBukhsh, A. M. Zoubir, and H. Koeppl, "Inverse reinforcement learning in swarm systems," in *Proc. of the 16th Conf.* on Auton. Agents and MultiAgent Sys., ser. AAMAS '17, 2017, p. 1413–1421.
- [20] E. Tolstaya, "Scalable learning in distributed robot teams," Ph.D. dissertation, University of Pennsylvania, 2021.
- [21] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang, "Physics-informed machine learning," *Nature Reviews Physics*, vol. 3, pp. 422–440, 06 2021.
- [22] R. Hasani, M. Lechner, A. Amini, D. Rus, and R. Grosu, "Liquid time-constant networks," *Proc. of the AAAI Conf. on Artificial Intell.*, vol. 35, no. 9, pp. 7657–7666, May 2021.
- [23] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, "Neural ordinary differential equations," in *NeurIPS*, 2018, pp. 6572–6583.
- [24] O. Khatib, "Real-time obstacle avoidance for manipulators and mobile robots," *The Int. J. of Robot. Res.*, vol. 5, no. 1, pp. 90–98, 1986.
- [25] A. Wikner, J. Pathak, B. Hunt, M. Girvan, T. Arcomano, I. Szunyogh, A. Pomerance, and E. Ott, "Combining machine learning with knowledge-based modeling for scalable forecasting and subgrid-scale closure of large, complex, spatiotemporal systems," *Chaos*, vol. 30, p. 053111, 05 2020.
- [26] SebLague, "Boids," https://github.com/SebLague/Boids/tree/master, 2019.
- [27] P. Holmes, J. L. Lumley, and G. Berkooz, *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*, ser. Cambridge Monographs on Mechanics. Cambridge University Press, 1996.
- [28] J. M. Joyce, Kullback-Leibler Divergence. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 720–722.