

Relative Distributed Formation and Obstacle Avoidance with Multi-agent Reinforcement Learning

Yuzi Yan[#], Xiaoxiang Li[#], Xinyou Qiu[#], Jiantao Qiu[†], Jian Wang[#], Yu Wang[#] and Yuan Shen[#]

Abstract—Multi-agent formation as well as obstacle avoidance is one of the most actively studied topics in the field of multi-agent systems. Although some classic controllers like model predictive control (MPC) and fuzzy control achieve a certain measure of success, most of them require precise global information which is not accessible in harsh environments. On the other hand, some reinforcement learning (RL) based approaches adopt the leader-follower structure to organize different agents' behaviors, which sacrifices the collaboration between agents thus suffering from bottlenecks in maneuverability and robustness. In this paper, we propose a distributed formation and obstacle avoidance method based on multi-agent reinforcement learning (MARL). Agents in our system only utilize local and relative information to make decisions and control themselves distributively, and will reorganize themselves into a new topology quickly in case that any of them is disconnected. Our method achieves better performance regarding formation error, formation convergence rate and on-par success rate of obstacle avoidance compared with baselines (both classic control methods and another RL-based method). The feasibility of our method is verified by both simulation and hardware implementation with Ackermann-steering vehicles.

I. INTRODUCTION

Formation control while avoiding obstacles is one of the most basic function of an multi-agent system (MAS). In scenarios like Internet of Vehicles, the autonomous platooning (as a formation task) and overtaking (as an obstacle avoidance task) are the most common and important maneuvers.

Most previous studies [1]–[6] regard the whole task as an optimization problem to plan the agents' route and movement according to the destination and reward function, while under constraints like avoiding obstacles and other agents during the motion. Since the optimization problem tends to be non-convex, some related works based on classic hierarchical control like *model predictive control (MPC)* [7], [8] or *fuzzy control* [9] are proposed to deal with the problem. However, most of them require high-precision global information like GPS and digital maps, leading to inapplicability in harsh environments such as search-and-rescue in emergency disasters. Besides, the collaboration between agents are not fully considered in these traditional methods, which means there is still a lot of room for improvement in terms of multi-agent collaboration.

*This work is supported by National Key R&D Program of China 2020YFC1511803 “Basic Research Strengthening Program of China (173 Program) (2020-JCJQ-ZD-015-01)” Tsinghua University - Meituan Joint Institute for Digital Life.

[#] Department of Electronic Engineering, Tsinghua University, Beijing, China. yyz21@mails.tsinghua.edu.cn, shenyuan.ee@tsinghua.edu.cn.

[†] This work is finished during his Ph.D career at Tsinghua University. He is now at Shanghai AI Laboratory, Shanghai, China.

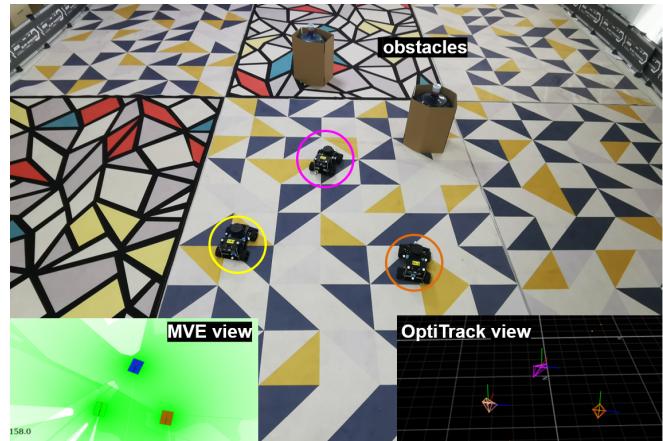


Fig. 1. The scenario of our problem. Ackermann-steering agents are placed randomly and meant to complete a formation task while avoiding obstacles detected by Lidar sensors. The sub-picture in the bottom left shows the results obtained by our self-developed simulator, called Multi-Vehicle Environment(MVE). The sub-picture in the lower right shows the precise coordinates of agents obtained by OptiTrack.

During the past few decades, the maturity of intelligent agent has been largely enhanced thanks to the deep combination of reinforcement learning (RL) and control theory [10]–[13]. Some previous works [14], [15] conduct RL to realize automatic formation control and obstacle avoidance. But most of them fail to get rid of the leader-follower structure, and mainly focus on controlling a certain single agent [16]–[20]. If the leader agent is destroyed or disconnected, the whole system will collapse. What's more, common RL-based works are verified only by numerical simulation. A few works implement their algorithms on hardware platforms but only take the omnidirectional wheel model into consideration [21], [22], which is not enough to be used in practical Ackermann-steering¹ [23] vehicular system.

In this paper, we design a distributed formation and obstacles avoidance algorithm with Multi-Agent Proximal Policy Optimization (MAPPO) [24], [25]. The proposed method requires only relative information which is easily available for real systems, and the control policy can be executed distributively. The **contributions** of our work are as follows: 1) Distributed. We put forward a relative formation strategy that is independent of global position information. Agents adjust their postures by taking into account the network

¹Ackermann-steering geometry is the practical model as human driven vehicles. Compared with other omnidirectional wheel systems, it has more constraints and faces challenges in convergence and time consumed for training.

topology obtained by spacial-temporal cooperation rather than absolute coordinates or orientation angle information. We avoid the leader-follower structure and train a policy that supports decentralized execution, which is more robust. 2) Adaptive. Multiple formation strategies are integrated in our model through policy distillation. If any agent in the formation is destroyed or disconnected, other agents will reorganize themselves into a new topology adaptively to continue their work. Besides, we use curriculum learning [26] to accelerate training regarding obstacle avoidance, improving the stability of the MAS in a complex environment. 3) Effective. We compare our method with existing traditional control algorithms (MPC [7], fuzzy control [9]) and a RL-based leader-follower method [19] by numerical simulations. The results of average formation error, formation convergence rate and success rate of obstacle avoidance illustrate that our method achieves better performance than the baselines. 4) Practical. We conduct our method on a hardware platform using intelligent cars with Ackermann-steering geometry as agents.

Notation: Throughout this paper, variables, vectors, and matrices are written as italic letters x , bold italic letters \mathbf{x} , and bold capital italic letters \mathbf{X} , respectively. Random variables and random vectors are written as sans serif letter x and bold letters \mathbf{x} , respectively. The notation $\mathbb{E}_{\mathbf{x}}\{\cdot\}$ is the expectation operator with respect to the random vector \mathbf{x} , and $\mathbb{1}(\cdot)$ is the indicator function which equals 1 if the condition is true and equals 0 otherwise.

II. PROBLEM FORMULATION

A. Relative Localization and Formation Error

High-precision location information is a prerequisite and important guarantee for complex tasks such as formation and obstacle avoidance. The state of the art studies mainly focus on the global localization optimization which is high-cost and unguaranteed in harsh environments [27]. In scenarios like Internet of Vehicles, people pay more attention to relative relationships which reflect the shape of network geometry [28], since the relative topology is sufficient to complete maneuvers like overtaking or formation and much easier to be obtained.

Considering a two-dimensional formation that consists of N agents and the set of agents is denoted as \mathcal{N} . The global position of agent i is denoted as $\mathbf{p}_i = [x_i \ y_i]^T$. In the local coordinate system of agent i , the relative position of any other agent j is denoted as $\mathbf{p}_{i \leftarrow j} = [x_j - x_i \ y_j - y_i]^T$. The relative position parameter vector of the formation is denoted as $\mathbf{p}_{i \leftarrow \tilde{i}} = [\mathbf{p}_{i \leftarrow 1} \dots \mathbf{p}_{i \leftarrow N}]^T$. For a given formation positioned as \mathbf{p} , the *equivalent geometry* of this formation is denoted as

$$T_{\omega}(\mathbf{p}) = (\mathbf{I}_N \otimes \Gamma_{\vartheta}) \mathbf{p} + \mathbf{1}_N \otimes \mathbf{t}, \quad (1)$$

where $\omega = [\Gamma_{\vartheta} \ \mathbf{t}]$, Γ_{ϑ} denotes the rotation matrix of angle $\vartheta \in [0, 2\pi]$, and $\mathbf{t} \in \mathbb{R}^2$ denotes the translation in x, y axes.

Since two formations are considered equivalent if one can be transformed into another through rigid body transformation like translation and rotation. The formation error $\mathcal{E}(\mathbf{p}, \mathbf{q})$

between two given formations \mathbf{p} and \mathbf{q} is defined as the squared Euclidean distance between the equivalent geometry sets $T_{\omega}(\mathbf{p})$ and $T_{\omega}(\mathbf{q})$:

$$\begin{aligned} \mathcal{E}(\mathbf{p}, \mathbf{q}) &= \min_{\mathbf{p} \in T_{\omega}(\mathbf{p}), \mathbf{q} \in T_{\omega}(\mathbf{q})} \|\mathbf{p} - \mathbf{q}\|^2 \\ &= \min_{\omega^*} \|\mathbf{p} - T_{\omega^*}(\mathbf{q})\|^2 \end{aligned} \quad (2)$$

Remark 1: It's obvious that $T_{\omega}(\mathbf{p}_{i \leftarrow \tilde{i}}) = T_{\omega}(\mathbf{p}_{j \leftarrow \tilde{j}})$ for any two agents i and j in the same formation topology. To keep the notation simple, we use $\check{\mathbf{p}}$ to represent the relative topology obtained by agent i , $\forall i \in \mathcal{N}$.

B. Multi-agent Reinforcement Learning

The relative distributed formation as well as obstacle avoidance can be regarded as a fully-cooperative problem, which is solved under a MARL framework.

The MARL process of N agents can be modeled as an extension of N Markov decision processes [29]. It composes of a state space \mathcal{S} describing the possible configurations of all agents, a set of actions $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_N$ and a set of observations $\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_N$.

Following [19], the action output of each agent i includes 4 control variable: $[\omega_L, \omega_R, v_F, v_B] \in \mathcal{A}_i$, where ω_L and ω_R represent the angular velocity of turning left and right; v_F and v_B represent the speed of moving forward and backward.

As for the observations, \mathcal{O}_i consists of variables as follow: 1) the relative distance d_{ij} and angle θ_{ij} towards other agent $j, j \in \mathcal{N}$, by which we estimate the network topology and calculate the formation error \mathcal{E} ; 2) the relative distance D_i and angle Θ_i towards the destination; 3) the shortest distance of nearby obstacles detected by the Lidar sensor $d_{im} = \min d_i$ and its corresponding direction θ_{im} according to the agent's coordinate, where $\mathbf{d}_i = [d_{i1} \ d_{i2} \ \dots \ d_{iM}]^T$ is the detected result with angle resolution $2\pi/M$.

Each agent i obtains their own reward $r_i(s^t, a_i^t) : \mathcal{S} \times \mathcal{A}_i \rightarrow \mathbb{R}$ to measure the feedback cost when taking action a_i^t with state s^t at time slot t , where a_i^t and s^t are always omitted for the simplicity of notion. A *policy*, denoted as $a_i \sim \pi_{\theta_i}(\cdot | s) : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A}_i)$, projects states to the probability measures on \mathcal{A}_i which returns the probability density of available state and action pairs (s, a_i) . θ_i refers to parameters of the function π , and is always omitted for the simplicity.

RL involves estimating the total expected reward $\eta(\pi) = \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r(s^t, a^t)]$ with policy gradient methods (e.g. PPO [24]), who optimizes the parameter θ of the policy according to the explored data.

Among all policy gradient methods, PPO shows its efficiency in both stabilizing the policy and exploring for optimal results. Moreover, it has become the most powerful baseline algorithm of DRL due to its generalization to various tasks, including MAS. In rest of the paper, our method applies MAPPO [25], an advanced version of PPO for multi-agent tasks to solve the formation problem.

III. APPROACH

A. Reward Function Design

One of the most important tasks of RL is the design of reward function. The proposed reward function is divided

into three parts: relative formation reward, navigation reward and obstacle avoidance reward

$$r = r_{\text{form}} + \alpha r_{\text{navi}} + \beta r_{\text{avoid}}, \quad (3)$$

where α and β are the hyper-parameters for reward trade-off.

The relative formation reward is designed based on the formation error as defined in (2). Given an ideal formation topology \mathbf{q} , the formation error $\mathcal{E}(\check{\mathbf{p}}, \mathbf{q})$ can be used to measure the difference between the ideal formation and the actual formation, which is an optimizable goal in the RL framework. In the proposed method, the relative formation reward is defined as:

$$r_{\text{form}} = -\frac{\mathcal{E}(\check{\mathbf{p}}, \mathbf{q})}{G(\mathbf{q})}. \quad (4)$$

$G(\mathbf{q}) = \max_{i,j \in \mathcal{N}} d_{ij}^2$ is a normalization factor related only to the size and topology of the ideal formation, where d_{ij} represents the distance between agent i and agent j in the formation. It is introduced to normalize the reward regardless of different formation scales, which avoid repeated adjustment of the hyper-parameters α and β .

The navigation reward is designed to reflect the efficiency of the agent moving towards the destination. Following [30], the navigation reward is designed as $r_{\text{navi}} = D_i^{t-1} - D_i^t$, where D_i^t represents the distance from agent i to the destination at time slot t .

The obstacle avoidance reward is designed to reflect the success rate of collision avoidance. Following [31], the number of collisions in time slot t is counted:

$$R_{\text{avoid}} = -\sum_{i \in \mathcal{N}, j \in \mathcal{M}} \mathbb{1}(d_{ij} < \delta_{ij}),$$

where δ_{ij} is the collision margin of agent i with entity j , where \mathcal{N} and \mathcal{M} denotes the set of agents and obstacles, respectively. Since Lidar or ultrasonic sensors can not distinguish whether the detected object is an agent or an obstacle, we treat all entities equally to be avoided.

B. Policy Distillation for Formation Adaptation

We call it *formation adaptation* that agents need to reorganize themselves into a new topology in case that any agent is disconnected. To achieve formation adaptation, we conduct policy distillation [32], a method that integrate the learned policies to handle different numbers of agents to complete the formation. For example, if one agent in a regular pentagonal formation completed by 5 agents is damaged due to collision of obstacles, the remaining 4 agents need to reorganize themselves into a square to move on.

The maximum number of agents in the formation is set to be N , and the corresponding ideal formation topology as \mathbf{q}_N . For n agents participating in the formation ($n < N$), we pre-set its corresponding ideal formation topology as \mathbf{q}_n . In order to allow different numbers of agents to complete the formation adaptively, we train teacher models according to different n and \mathbf{q}_n . Then we use policy distillation to train

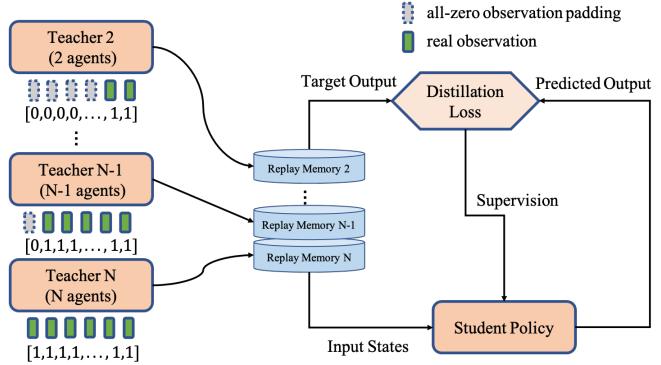


Fig. 2. Policy distillation for formation adaptation. The observations (grey blocks for all-zero observation padding and green blocks for real observation) as well as corresponding death masking are input into teacher models for training. Teacher models generate input states and target output for supervised learning of student model where KL divergence is used as loss function.

a student model that can handle multiple situations from the teacher models.

We train several teacher models which share the same structure but with different n and \mathbf{q}_n . Connected agents are agents participating in the formation. We set *death masking* [25] to mark the status of agents. The corresponding value of a connected agent in the *death masking* is set to be 1 while that of a disconnected agent is set to be 0. Due to the different number of agents participating in the formation, the length of the input observations (such as the relative distance and relative angle of other agents) is different for different teacher models. To align the length of observation inputs, all-zero observation padding caused by *death masking* is set for disconnected agents as shown in Fig. 2.

In the stage of policy distillation, the trained teacher models produce inputs and targets, which are then stored in separate memory buffers. The student model learns from the data stored sequentially, switching to a different memory buffer every episode, just as in [33]. We adopt the distillation setup of [32] and use the KL divergence as the loss function:

$$L(p_T(s)|p_S(s)) = \sum_{i=1}^{|\mathcal{A}|} p_{k,T}(s) \ln \frac{p_{k,T}(s)}{p_{k,S}(s)},$$

where $|\mathcal{A}|$ is the dimension of the action space, $p_{k,S}(s)$ and $p_{k,T}(s)$ represent the probabilities for action k of the student model S and the teacher model T given state observation s , respectively.

C. Curriculum Learning for Obstacle Avoidance

As mentioned in Section II, the optimization problem tends to be non-convex. Training the network directly in a complex environment is likely to cause non-convergence. In order to reduce the difficulty of training and speed up convergence, we use the idea of *predefined curriculum learning* [34] and set different levels of difficulties for obstacle avoidance.

The specific curriculum settings are as follows: We set a total of 5 difficulty levels according to the density of

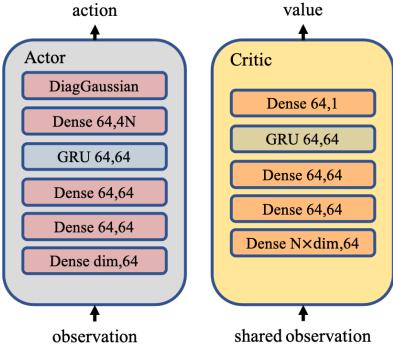


Fig. 3. The model structure of actor-critic networks we use. Note that we also use tricks such as observation normalization, layer normalization, ReLU activation, etc., just as in [25], but do not show here mainly for simplicity.

obstacles. At the beginning, no obstacle will be placed in the environment (level-0). The agents focus mainly on learning basic formation strategies and navigating themselves to the destination. As the reward curve gradually converges, the level of difficulty will increase and more obstacles of different sizes will gradually appear in the environment. The agents are able to learn how to avoid obstacles to reach the destination as well as maintaining the ideal formation as stably as possible. The settings and results will be introduced in detail in the next section.

IV. EXPERIMENTS AND RESULTS

To validate the effectiveness of our algorithm, we compare our method with several baseline methods including MPC [7], fuzzy control [9], and a RL-based leader-follower scheme [19] on the Multi-agent Particle-world Environment (MPE) [29]. To further evaluate the practicability of our method under physical constraints, we not only develop a new simulator, called Multi-Vehicle Environment (MVE)², which supports the Ackermann-steering model rather than omnidirectional wheel model, but also implement our algorithm on a corresponding hardware platform.

A. Model Configuration

In our model, the actor network and the critic network both consist of 3 dense layers, 1 GRU layer and 1 dense layer in order, as shown in Fig. 3. The hidden sizes of the dense layer and GRU layer are set to be 64. We follow the common practices in PPO implementation, including Generalized Advantage Estimation (GAE) [35] with advantage normalization, observation normalization, gradient clipping, layer normalization, ReLU activation with orthogonal initialization. Following the PopArt technique proposed by [36], we normalize the values by a running average over the value estimates to stabilize value learning.

Following the centralized training decentralized execution principle [37], the shared observation (i.e. observations from all agents) is input to critic network. In training process, we

take 30,000,000 steps to optimize our model in total. We train our model on 1 NVIDIA RTX3090 GPU and AMD EPYC 7R32 48-Core CPU. The Adam optimizer [38] is used with $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-3}$.

Remark 2: The dynamics design of the vehicle takes into account the characteristics of the Ackerman-steering model. In the MPE environment, the model is simplified by setting wheelbase to 0. In the MVE environment, we designed the vehicular body based on the actual hardware, which will be described in detail at Section IV-C.

B. Simulation Results in MPE

The scenario settings are as follows: The maximum speed of agents is restricted to be 1m/s. Agents are squares with a side length of 0.01m. Agents are randomly initiated in the range of $x \in [-2m, 2m]$, $y \in [-2m, 2m]$ and the destination will be initiated at a distance of at least 36m from the agent's initial location. Obstacles are circles with random radius between [0.01m, 0.05m]. Obstacles are placed with uniform distribution between the start point [0m, 0m] and the destination. The amount of the obstacles randomly distributed on the map will increase according to different difficulty levels. The rendering of our scenario in MPE is shown in Fig. 5.

1) *Training Process:* As is shown in Fig. 4, we train the corresponding models for the formation with 3, 4 and 5 agents, respectively, where the ideal formation topology is set to be regular polygon. It is demonstrated that in the later stage of training, our model is able to organize and maintain formation finely (the formation error rises close to 0) while avoiding the obstacles and advancing to the destination.

2) *Baseline Comparison:* We compare our proposed method with other methods, including: 1) MPC [7], 2) fuzzy control [9], 3) RL-based leader-follower scheme [19]. In baseline comparison, the agent number is set to be 4, and the difficulty level is set to be level-3. The ideal formation topology is set to be a square with side length of 4m.

The traditional controllers (MPC and fuzzy control) treat the multi-task scenario as a motion-planning problem. The desired position of each sub-task will be taken as input to the controller. To ensure the performance of compound behaviors, the sub-tasks are assigned to different priorities with fine-tuned thresholds. The controller will execute the task with the highest priority at a time step. Since safety is the first concern, the avoiding behavior is the top priority among all, while formation behavior comes second. Thus, the navigation behavior won't be triggered unless the agents enter a safety area and the formation error converges. The obstacle-avoidance behavior is designed with a stream-based path planner [39].

The RL-based leader-follower scheme uses the same network architecture as our method to train the agents to form the ideal topology by tracking a virtual leader. Each agent is only responsible of tracking its relative position towards the leader and obstacles, as demonstrated in [19].

Fig. 6 shows the convergence of the formation error and the navigation reward executed by different controllers. It

²<https://github.com/efc-robot/MultiVehicleEnv>.

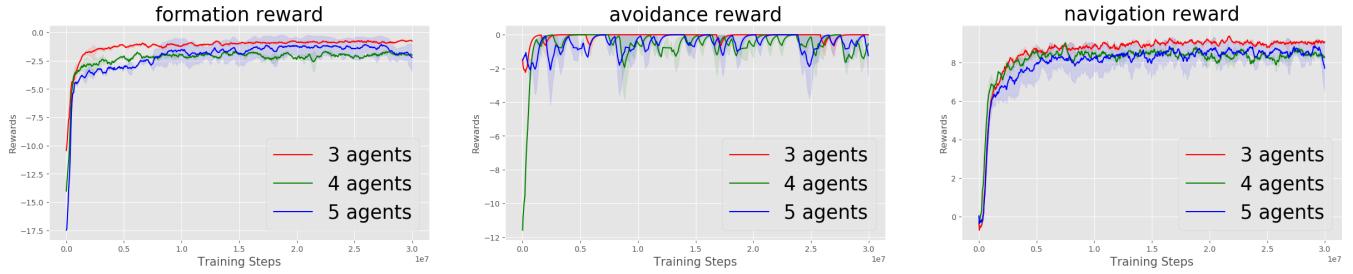


Fig. 4. The curve of the reward function (formation, obstacle avoidance and navigation) when different numbers of agents are in formation. Due to normalization, the trade-off parameters α and β defined in (3) do not need to be adjusted manually as the number of agents changes.

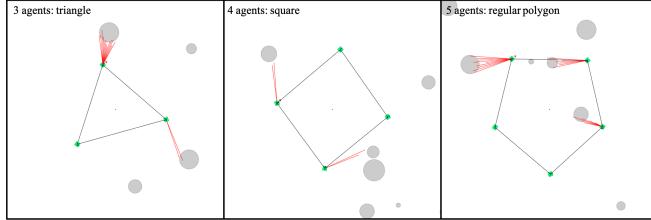


Fig. 5. Simulation results in Multi-agent Particle-world Environment (MPE). The agents (green squares) are navigating towards the destination that is out of the screen, as well as maintaining the formation and avoiding obstacles (grey circles). The red rays represent the field of view of the equipped Lidar.

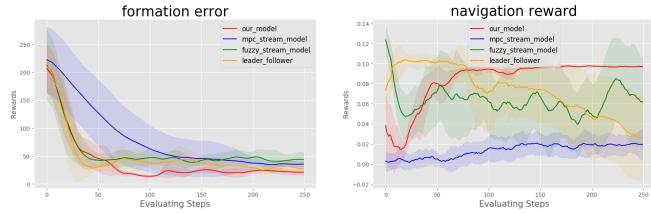
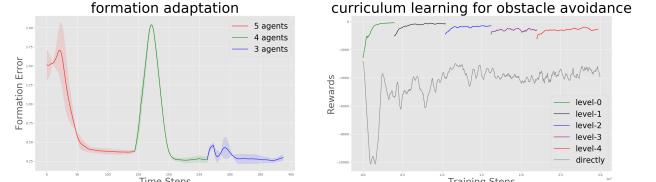


Fig. 6. Baseline comparison. Compared with MPC (blue), fuzzy control (green) and leader-follower method (yellow), our model (red) achieves the highest formation accuracy and stability (left side). Meanwhile, our method maintains the most stable and highest navigation speed without any agent crashing obstacles (right side).

appears that our method outperforms in formation accuracy, convergence speed and ability to handle multi-tasks. Specifically, our method can converge to a lower formation error, which indicates better maintenance of formation during movements. The navigation rewards depict the advantage of DRL scheme when dealing with multi-task problems. The formation is encouraged to accelerate the navigation as the formation error converges, yet the traditional controllers struggle to switch between different behaviors. It is mainly because our method jointly optimizes all reward functions, which enables the agents to successfully learn the priority between formation control and navigation without any prior knowledge. Our method also achieves an on-par success rate of obstacle avoidance with the baselines and achieve higher efficiency of moving to the destinations at the same time.

3) Formation Adaptation: The formation policies based on a fixed agent number can be used as teacher models to train a student policy, which is supposed to guide the agents to be reorganized into new formations when some agents are



(a) Change of formation error over time steps in formation adaptation. In step 144 and 263, two agents gray line represents the reward curve are disconnected in sequence. The ideal formation is changed from colorful lines represent the reward regular pentagon to square to tri-angle for training directly in level-4. The angle accordingly. (b) Promotion of training performance by curriculum learning. The difficulty level is changed from level-0 to level-4 in a step-wise process.

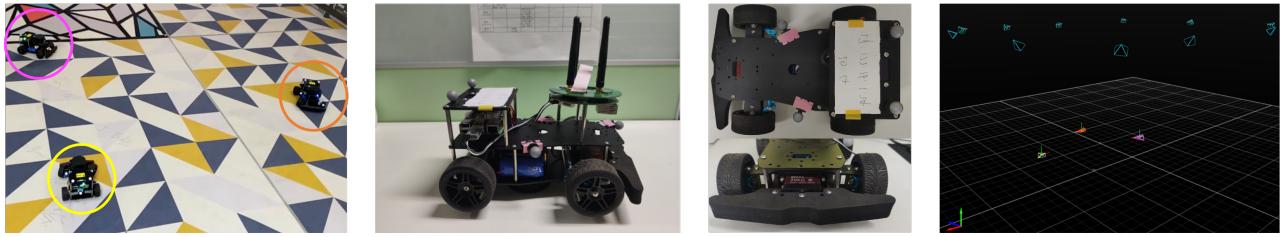
Fig. 7. Results of formation adaptation and curriculum learning.

disconnected accidentally.

In practice, we set the maximum number of agents in the formation to be 5, and use the formation policies with 3, 4, and 5 agents as the teacher models, respectively. α and β in (3) are set to be 5 and 10. The student model shares the same structure with the teacher model. The sample batch size is set to be 1000. The Adam optimizer is used with $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-3}$. The student model is trained for 5000 episodes to reach the convergence. The KL divergence is used as loss metric.

We conduct experiments to test the ability of the student policy to guide agents to reorganize their formation when some of them are disconnected. As is shown in Fig. 7(a), at the beginning of the experiment, 5 agents are guided to form a pentagon formation while moving forward to the destination. In step 144 and 263, we deliberately disconnect 2 agents in sequence, and the formation will change to square and triangle automatically. Due to changes of the ideal formation topology, the formation error will fluctuate, but new formation will be reorganized and the formation will converge quickly.

4) Curriculum Learning for Obstacle Avoidance: Following [26], we set different difficulty levels for obstacle avoidance. 5 levels are set in total from level-0 to level-4, and the corresponding obstacle density is $0/m^2$, $1 \times 10^{-2}/m^2$, $2 \times 10^{-2}/m^2$, $3 \times 10^{-2}/m^2$, and $5 \times 10^{-2}/m^2$. Initially, no obstacle will be generated in the environment. When the model gets a converged reward under the current level, the difficulty of the task will be enhanced. More obstacles will



(a) The formation scenario of 3 intelligent vehicle. (b) The overview of the intelligent vehicle equipping with Lidar. (c) The chassis and (d) The OptiTrack motion capture system used for getting the position groundtruth.

Fig. 8. The hardware platform that is fully consistent with MVE.

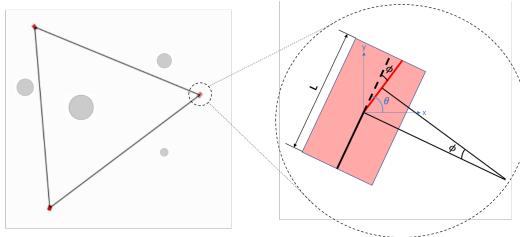


Fig. 9. Simulation results in MVE. Left side: The agents (pink rectangles) are forming into a triangle as well as avoiding obstacles (grey circles). Right side: The schematic diagram of Ackermann-steering geometry.

be generated in the scene to improve the obstacle avoidance ability of the agent cluster. It takes around 6,000,000 steps to train each difficulty level.

As is shown in Fig. 7(b), it is found that if we directly train the policy for formation and obstacle avoidance in complex scenarios with too many obstacles, the model will converge slower due to poor exploration, and the final average reward will be much lower. By curriculum learning, agents can achieve better formation and obstacle avoidance performance in the final level (level-4).

5) Asymmetric Formation: In addition to formations with regular shapes (e.g. triangle, square, regular polygon), our algorithm can also complete formations with asymmetric topology (e.g. irregular convex polygon). It is found that separated policy is better than shared policy since the relative position of each node is not completely symmetrical.

C. Hardware Implementation with MVE

The agent in MVE is designed to be an intelligent vehicle with Ackermann-steering. Different from the simplified model in MPE, we consider the constraints of wheelbase and steering angle in practice, which leads to the nonzero turning radius of the vehicle. Given the wheelbase L and the steering angle ϕ , the turning radius R can be calculated as: $R = \frac{L}{\tan \phi}$. Given the speed of the rear wheels v_b , the angular velocity of the orientation angle θ can be calculated as $\omega = \frac{v_b}{R} \tan \phi$. If L is set to be 0, model will degrade to MPE and we can directly control the vehicle with v and ω . In MVE and hardware implementation, the control variables are changed to θ and v_b . Due to constraints of the hardware control system, we discretize the value of the control variable. The decision frequency is 1Hz while

the control frequency is 100Hz. At the decision-making stage, the model gives the ideal v_b and ϕ . At the control stage, the agent will try its best to reach the ideal control variables. The detailed settings can be found at <https://github.com/efc-robot/MultiVehicleEnv>.

In hardware implementation, we use scenario as Fig. 1 and Fig. 8(a). Multiple Ackermann-steering agents are placed in a room randomly. Cylindrical objects with a radius of 0.14m are placed in the field as obstacles. The intelligent vehicles are as Fig. 8(b). The wheelbase L is 0.20m, the overall width is 0.18m and the overall length is 0.25m, which is fully consistent with the simulation platform as is shown in Fig. 8(c). The max speed is constrained to 0.361m/s and the max steering angle is constrained to 0.298rad/s. We use an OptiTrack motion capture system³ to get the groundtruth of positions and give relative observations as Fig. 8(d). Agents detect obstacles by Lidars. We conduct several hardware experiments with different numbers of agents forming in an environment with obstacles avoidance. Demos can be found at <https://sgroupresearch.github.io/relativeformation/>.

V. CONCLUSION

In this paper, we develop a MAPPO-based distributed formation and obstacle avoidance method, in which agents only use their local and relative information to make movement decisions. We introduce policy distillation to make the formation system adaptive in case of agents' accidental disconnection. Curriculum learning is also used to simplify the learning process. Our model achieves better performance regarding average formation error, formation convergence rate and success rate of obstacle avoidance. Besides, we also build a new simulation environment and a supporting hardware platform with Ackermann-steering geometry to verify the feasibility of our algorithm.

For the future work, we will explore large-scale distributed formation methods where agents are not fully connected and can only get the information with the neighbors. Besides, we will also concentrate on developing our self-developed MVE and corresponding hardware platform in order to solve the sim-to-real problems in MARL algorithm deployments.

³<https://www.optitrack.com>.

REFERENCES

- [1] K.-K. Oh, M.-C. Park, and H.-S. Ahn, "A survey of multi-agent formation control," *Automatica*, vol. 53, pp. 424–440, 2015.
- [2] Z. Sun, M.-C. Park, B. D. Anderson, and H.-S. Ahn, "Distributed stabilization control of rigid formations with prescribed orientation," *Automatica*, vol. 78, pp. 250–257, 2017.
- [3] Z. Lin, W. Ding, G. Yan, C. Yu, and A. Giua, "Leader-follower formation via complex laplacian," *Automatica*, vol. 49, no. 6, pp. 1900–1906, 2013.
- [4] P. Ogren and N. E. Leonard, "Obstacle avoidance in formation," in *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*, vol. 2. IEEE, 2003, pp. 2492–2497.
- [5] C. De La Cruz and R. Carelli, "Dynamic model based formation control and obstacle avoidance of multi-robot systems," *Robotica*, vol. 26, no. 3, pp. 345–356, 2008.
- [6] G. Wen, C. P. Chen, and Y.-J. Liu, "Formation control with obstacle avoidance for a class of stochastic multiagent systems," *IEEE Transactions on Industrial Electronics*, vol. 65, no. 7, pp. 5847–5855, 2017.
- [7] X. Li, K. Ma, J. Wang, and Y. Shen, "An integrated design of cooperative localization and motion control," in *Proc. IEEE Int. Conf. Commun. in China*, Changchun, China, Aug. 2019, pp. 1–5.
- [8] L. Dai, Q. Cao, Y. Xia, and Y. Gao, "Distributed mpc for formation of multi-agent systems with collision avoidance and obstacle avoidance," *Journal of the Franklin Institute*, vol. 354, no. 4, pp. 2068–2085, 2017.
- [9] J. E. Naranjo, C. Gonzalez, R. Garcia, and T. De Pedro, "Lane-change fuzzy control in autonomous vehicles for the overtaking maneuver," *IEEE Trans. Intell. Transp. Syst.*, vol. 9, no. 3, pp. 438–450, 2008.
- [10] H. Iima and Y. Kuroe, "Swarm reinforcement learning methods improving certainty of learning for a multi-robot formation problem," in *2015 IEEE Congress on Evolutionary Computation (CEC)*, 2015, pp. 3026–3033.
- [11] Z. Sui, Z. Pu, J. Yi, and S. Wu, "Formation control with collision avoidance through deep reinforcement learning using model-guided demonstration," *IEEE Trans. Neural Netw. Learning Syst.*, vol. 32, no. 6, pp. 2358–2372, 2020.
- [12] F. Kobayashi, N. Tomita, and F. Kojima, "Re-formation of mobile robots using genetic algorithm and reinforcement learning," in *SICE 2003 Annual Conference (IEEE Cat. No. 03TH8734)*, vol. 3. IEEE, 2003, pp. 2902–2907.
- [13] Z. Sui, Z. Pu, J. Yi, and X. Tan, "Path planning of multiagent constrained formation through deep reinforcement learning," in *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018, pp. 1–8.
- [14] G. Wen, C. P. Chen, J. Feng, and N. Zhou, "Optimized multi-agent formation control based on an identifier-actor-critic reinforcement learning algorithm," *IEEE Trans. Fuzzy Syst.*, vol. 26, no. 5, pp. 2719–2731, Oct. 2018.
- [15] M. Knopp, C. Aykin, J. Feldmaier, and H. Shen, "Formation control using gq (λ) reinforcement learning," in *26th IEEE Int. Symp. Robot and Human Interactive Commun. (RO-MAN)*. IEEE, 2017, pp. 1043–1048.
- [16] Y. Zhou, F. Lu, G. Pu, X. Ma, R. Sun, H.-Y. Chen, and X. Li, "Adaptive leader-follower formation control and obstacle avoidance via deep reinforcement learning," in *2019 IEEE/RSJ Int. Conf. Intell. Robots and Syst. (IROS)*. IEEE, 2019, pp. 4273–4280.
- [17] E. Yang and D. Gu, "A multiagent fuzzy policy reinforcement learning algorithm with application to leader-follower robotic systems," in *2006 IEEE/RSJ Int. Conf. Intell. Robots and Syst.* IEEE, 2006, pp. 3197–3202.
- [18] M. S. Miah, A. Elhussein, F. Keshtkar, and M. Abouheaf, "Model-free reinforcement learning approach for leader-follower formation using nonholonomic mobile robots," in *The 33rd Int. Flairs Conf.*, 2020.
- [19] X. Qiu, X. Li, J. Wang, Y. Wang, and Y. Shen, "A drl based distributed formation control scheme with stream based collision avoidance," *arXiv preprint arXiv:2109.03037*, (2021).
- [20] Z. Han, K. Guo, L. Xie, and Z. Lin, "Integrated Relative Localization and Leader-Follower Formation Control," *IEEE Trans. Autom. Control*, vol. 64, no. 1, pp. 20–34, Jan. 2019.
- [21] C. Ren and S. Ma, "A continuous dynamic model for an omnidirectional mobile robot," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May/Jun 2014, pp. 2919–2924.
- [22] G. Indiveri, "Swedish wheeled omnidirectional mobile robots: Kinematics analysis and control," *IEEE Trans. Robot.*, vol. 25, no. 1, pp. 164–171, 2009.
- [23] W. C. Mitchell, A. Staniforth, and I. Scott, "Analysis of ackermann steering geometry," SAE Technical Paper, Tech. Rep., 2006.
- [24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [25] C. Yu, A. Velu, E. Vinitksy, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of mappo in cooperative, multi-agent games," *arXiv preprint arXiv:2103.01955*, 2021.
- [26] S. Narvekar, B. Peng, M. Leonetti, J. Sinapov, M. E. Taylor, and P. Stone, "Curriculum learning for reinforcement learning domains: A framework and survey," *arXiv preprint arXiv:2003.04960*, 2020.
- [27] Y. Shen, H. Wyneersch, and M. Z. Win, "Fundamental limits of wideband localization – Part II: Cooperative networks," *IEEE Trans. Inf. Theory*, vol. 56, no. 10, pp. 4981–5000, Oct. 2010.
- [28] Y. Liu, Y. Wang, X. Shen, J. Wang, and Y. Shen, "UAV-aided relative localization of terminals," *IEEE J. Internet Things*, 2021, Early Access.
- [29] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Proc. Adv. Neural Inf. Process. Syst.*, pp. 6379–6390, 2017.
- [30] C. Wang, J. Wang, Y. Shen, and X. Zhang, "Autonomous navigation of UAVs in large-scale complex environments: A deep reinforcement learning approach," *IEEE Trans. Veh. Technol.*, vol. 68, no. 3, pp. 2124–2136, Mar. 2019.
- [31] Y. Jin, Y. Zhang, J. Yuan, and X. Zhang, "Efficient multi-agent cooperative navigation in unknown environments with interlaced deep reinforcement learning," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 2897–2901.
- [32] S. Green, C. M. Vineyard, and C. K. Koç, "Distillation strategies for proximal policy optimization," *arXiv preprint arXiv:1901.08128*, 2019.
- [33] A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirkpatrick, R. Pascanu, V. Mnih, K. Kavukcuoglu, and R. Hadsell, "Policy distillation," *arXiv preprint arXiv:1511.06295*, 2015.
- [34] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. 26th annual intern. conf. machine learning*, 2009, pp. 41–48.
- [35] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *arXiv preprint arXiv:1506.02438*, 2015.
- [36] M. Hessel, H. Soyer, L. Espeholt, W. Czarnecki, S. Schmitt, and H. van Hasselt, "Multi-task deep reinforcement learning with popart," in *Proc. AAAI Conf. Artif. Intell.*, vol. 33, no. 01, 2019, pp. 3796–3803.
- [37] G. Chen, "A new framework for multi-agent reinforcement learning—centralized training and exploration with decentralized execution via policy distillation," in *Proc. 19th Intern. Conf. Auto. Agents and MultiAgent Syst.*, 2020, pp. 1801–1803.
- [38] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [39] QiangWang, J. Chen, H. Fang, and Q. Ma, "Flocking control for multi-agent systems with stream-based obstacle avoidance," *Trans. Institute of Measurement and Control*, vol. 36, p. 391–398, 2014.