

# DL-DRL: A double-layer deep reinforcement learning approach for large-scale task scheduling of multi-UAV

Xiao Mao, Guohua Wu, and Mingfeng Fan

**Abstract**—This paper studies deep reinforcement learning (DRL) for the task scheduling problem of multiple unmanned aerial vehicles (UAVs). Current approaches generally use exact and heuristic algorithms to solve the problem, while the computation time rapidly increases as the task scale grows and heuristic rules need manual design. As a self-learning method, DRL can obtain a high-quality solution quickly without hand-engineered rules. However, the huge decision space makes the training of DRL models becomes unstable in situations with large-scale tasks. In this work, to address the large-scale problem, we develop a divide and conquer-based framework (DCF) to decouple the original problem into a task allocation and a UAV route planning subproblems, which are solved in the upper and lower layers, respectively. Based on DCF, a double-layer deep reinforcement learning approach (DL-DRL) is proposed, where an upper-layer DRL model is designed to allocate tasks to appropriate UAVs and a lower-layer DRL model [i.e., the widely used attention model (AM)] is applied to generate viable UAV routes. Since the upper-layer model determines the input data distribution of the lower-layer model, and its reward is calculated via the lower-layer model during training, we develop an interactive training strategy (ITS), where the whole training process consists of pre-training, intensive training, and alternate training processes. Experimental results show that our DL-DRL outperforms mainstream learning-based and most traditional methods, and is competitive with the state-of-the-art heuristic method [i.e., OR-Tools], especially on large-scale problems. The great generalizability of DL-DRL is also verified by testing the model learned for a problem size to larger ones. Furthermore, an ablation study demonstrates that our ITS can reach a compromise between the model performance and training duration.

**Index Terms**—Deep reinforcement learning, divide and conquer-based framework, interactive training, multi-UAV task scheduling.

## I. INTRODUCTION

Unmanned aerial vehicles (UAVs) are widely utilized in the domains of package delivery [1], environment surveillance [2], target tracking [3], and task reconnaissance [4] due to their high flexibility, strong mobility, and low power consumption. With the increasing task quantity and environment complexity, the multi-UAV system attracts great attention and gradually replaces the single UAV in mission execution. As the basis of UAVs cooperation in a multi-UAV system, task scheduling plays a significant role in ensuring UAVs complete tasks safely and effectively.

As a variant of the multiple traveling salesman problem (M-TSP) [5], the multi-UAV task scheduling problem is known to be NP-hard [6], which is difficult to find an optimal solution in a reasonable computation time. To overcome this challenge, previous studies have proposed exact and heuristic algorithms to solve this problem. Exact algorithms [7], [8], [9] can obtain the optimal solution of small-scale problems, while they are not suitable for large-scale problems due to the exponentially increasing computation time. Heuristic algorithms, such as genetic algorithm (GA) [10], ant colony optimization (ACO) [11], particle swarm optimization (PSO) [12], simulated annealing algorithm (SA) [13], etc. are desirable to address large-scale problems, which strikes a balance between solution quality and time efficiency. However, the heuristic rules need to be manually designed using prior knowledge, it could be nontrivial to design an effective algorithm for a specific problem. In addition, previous methods usually solve the problem in a whole manner, i.e., directly generating a task scheduling scheme that determines the execution order of tasks for each UAV, which encounters the phenomenon of the “curse of dimensionality” when solving large-scale problems.

Recently, to solve the large-scale task scheduling problem of multi-UAV, several studies decompose the original problem into some subproblems, which are then solved via conventional heuristic methods [14], [15]. In this way, the problem complexity can be effectively reduced, making these methods have the potential to solve large-scale problems. However, hand-engineered rules are still the main issue that limits the performance and generality of these methods. With the prosperity of deep learning (DL) and reinforcement learning (RL), deep reinforcement learning (DRL) is widely used in games [16], [17], robotics [18], and natural language processing [19]. Recently, DRL is also popularly applied in combinatorial optimization problems (COP), such as vehicle routing problem (VRP) [20], traveling salesman problem (TSP) [21], and orienteering problem (OP) [22]. DRL can self-learn a decision-making policy to quickly engender a high-quality solution without human intervention. However, when the problem scale increases, the decision space expands sharply, which may lead to an unstable training process of the DRL model, and make the model difficult to converge [23].

To address the aforementioned issues, in this paper, we design a divide and conquer-based framework (DCF) to decompose the original multi-UAV scheduling problem into a task allocation and a UAV route planning subproblems which are solved in the upper and lower layers, respectively. Based

Xiao Mao, Guohua Wu (*the corresponding author*) and Mingfeng Fan are with the School of Traffic and Transportation Engineering, Central South University, Changsha 410075, China. (e-mail: maoxiao@csu.edu.cn, guohuawu@csu.edu.cn, mingfan@csu.edu.cn)

on DCF, the decision space of each subproblem is considerably shrunk compared to the original problem, and we propose a double-layer deep reinforcement learning approach (DL-DRL) to solve the subproblems. DL-DRL contains two different DRL models (i.e., the upper-layer and lower-layer DRL models) which are constructed based on an encoder-decoder architecture [24]. The upper-layer model is constructed for solving the task allocation subproblem, in which an encoder based on the self-attention mechanism [25] and a UAV selection decoder are designed. The state-of-the-art attention model (AM) [26] is used in the lower layer to address the UAV route planning subproblem. Unlike hierarchical DRL which needs to learn the goal/option determination policy by trial and error [27], [28], the goals of our DL-DRL are determined by DCF in advance. In this way, the difficulty of policy learning is decreased as only the action selection policy needs to be learned. Furthermore, in our DL-DRL, the upper-layer and lower-layer models interact with each other, i.e., during the models' training process, the upper-layer model determines the input data distribution of the lower-layer model, and the reward of the upper-layer model needs to be calculated by the lower-layer model. Therefore, we firstly propose an interactive training strategy (ITS) for the efficient model training, in which the whole training process consists of pre-training, intensive training, and alternate training processes. Experiment results show the effectiveness and great generalization performance of our approach in the large-scale multi-UAV task scheduling problem. The main contributions of this paper can be summarized as follows.

- A divide and conquer-based framework (DCF) is designed to divide the multi-UAV task scheduling problem into the task allocation and UAV route planning subproblems. Based on DCF, we propose a DL-DRL, in which the upper-layer and lower-layer DRL models are developed and applied to solve different subproblems, respectively.
- An interactive training strategy (ITS) is constructed for the effective training of the upper-layer and lower-layer DRL models. This strategy constructs three training processes, including pre-training, intensive training, and alternate training, which efficiently train policy networks of the DRL models.
- Extensive experiment results indicate our approach outperforms learning-based baselines while enjoying high running speed compared to traditional exact and heuristic baselines. The superiority of our training strategy is also verified via an ablation study that it can reach a compromise between the model performance and training duration. Besides, our approach shows improved generalization performance in problems with larger-scale tasks.

The remainder of this paper is organized as follows. Section II briefly reviews the related research work on the multi-UAV task scheduling problem. Section III formulates the mathematical model and constructs the DCF for problem decomposition. Section IV constructs the Markov decision process (MDP) for the multi-UAV task scheduling problem and elaborates our DL-DRL approach. Section V conducts extensive experiments and results analysis. Section VI concludes the paper.

## II. RELATED WORK

As a variant of M-TSP, the multi-UAV task scheduling problem is a kind of COP which needs to be solved under some constraints like traveling distance, time window, platform capacity, etc. Researchers typically construct a mixed-integer linear program model [7] for this problem and design exact and heuristic algorithms to address the problem.

Exact algorithms such as branch and bound [29], branch and price [30], column generation [31], and dynamic programming [8] can obtain the optimal solution, while they are infeasible for solving large-scale problems due to the prohibitive computation time for exhaustive search. Heuristic algorithms that can reduce search space by carefully designed heuristic rules become popular alternatives. Ye et al. [32] developed the GA with a multi-type gene chromosome encoding scheme and an adaptive operation for efficient solution searching. Shang et al. [33] combined GA and ACO to propose a hybrid algorithm (i.e., GA-ACO), in which the bad individuals of the population in GA are replaced by the better ones in ACO during the population evolution. Chen et al. [34] designed a transposition and extension operation and then proposed a modified two-part wolf pack search (MTWPS) algorithm for solving the problem. However, the human-made heuristic rules deeply rely on the human experience and domain knowledge, hence the solution quality may not be guaranteed.

DRL which combines the advantages of DL and RL is widely utilized in multiple domains, including games, robotics, speech recognition, and natural language processing [35] due to its powerful learning capability. As DRL can learn a decision-making policy via the interaction with the environment, researchers have recently applied it to solve COP [36], [37], which avoids the manual design of heuristic rules. Bello et al. [38] constructed a DRL model based on pointer networks (PtrNet) [39] to solve the TSP, which breaks through the barriers of DRL to solve COP. To improve the performance of PtrNet-based DRL, Ma et al. [40] proposed a graph pointer network that combines the graph neural network and PtrNet for better task feature extraction. In addition, Kool et al. [26] proposed the AM method based on the Transformer architecture [25] to solve multi-type COP, which performs better than a wide range of learning-based and conventional baselines. Xu et al. [41] developed the AM to efficiently solve VRPs by a multiple relational attention mechanism. These DRL methods have a good performance and a short computation time in solving COP, while the training process becomes unstable in large-scale situations due to the huge decision space [23].

In summary, due to the explosively increasing computational complexity and hand-engineered rules, conventional exact and heuristic algorithms are frustrating in solving the large-scale task scheduling problem of multi-UAV. DRL which learns a decision-making policy by the interaction with the environment can be a good alternative. However, the explosive expansion of decision space leads to the instability of the DRL model training on large-scale problems. To tackle this issue, recent work has concentrated on the problem decomposition. Hu et al. [42] divided the M-TSP into several small-scale

TSP via the proposed DRL policy networks. Liu et al. [43] constructed a decomposition framework named EA-DRL to decompose the OP into a knapsack problem and a TSP which are solved by an evolutionary algorithm and a DRL, respectively. In this way, the large-scale problem is divided into several subproblems that DRL can solve efficiently due to their considerably smaller decision spaces. Motivated by this observation, we design a DCF that decomposes the large-scale multi-UAV task scheduling problem into two subproblems: task allocation and UAV route planning. Based on DCF, we propose a DL-DRL approach to solve these two subproblems using two different DRL models, which has huge potential for solving large-scale multi-UAV task scheduling problems.

### III. PROBLEM FORMULATION AND FRAMEWORK DESCRIPTION

In this section, a mixed-integer linear program model is developed for the multi-UAV task scheduling problem, where the objective is to maximize the total number of executed tasks of UAVs. To solve the large-scale problem, we decompose the original problem into a task allocation and a UAV route planning subproblems based on the divide and conquer principle, and then construct a DCF. The DCF consists of the upper and lower layers in which the two subproblems are solved, respectively.

#### A. Mathematical Model

Let  $T = \{0, 1, \dots, N\}$  be the task set, where 0 and  $N$  denote the UAV depot, UAVs must depart from the depot to execute tasks and return to it.  $U = \{1, 2, \dots, V\}$  is the UAV set, and  $V$  is the number of UAVs. In this paper, all UAVs are considered homogeneous, and three types of decision variables are defined in the mathematical model of the multi-UAV task scheduling problem: 1)  $x_{ijk}$  is a binary variable that equals one if the UAV  $k$  executes the task  $j$  right after the task  $i$  and zero otherwise; 2)  $y_{ik}$  is also a binary variable that equals one if the UAV  $k$  has executed the task  $i$  and zero otherwise; 3)  $z_{ik}$  is a continuous variable indicating the remaining flight distance of the UAV  $k$  when visiting task  $i$ . The mathematical model can be formulated as follows.

$$\max \sum_{i=1}^{N-1} \sum_{k=1}^V y_{ik} \quad (1)$$

Subject to:

$$\sum_{k=1}^V y_{ik} \leq 1, i \in \{1, 2, \dots, N-1\} \quad (2)$$

$$\sum_{j=1}^N x_{ijk} = y_{ik}, i \in \{1, 2, \dots, N-1\}, k \in U \quad (3)$$

$$\sum_{j=1}^N x_{ijk} = \sum_{j=0}^{N-1} x_{jik} \leq 1, i \in \{0, 1, \dots, N\}, k \in U \quad (4)$$

$$\sum_{k=1}^V \sum_{j=1}^{N-1} x_{0jk} = \sum_{k=1}^V \sum_{i=1}^{N-1} x_{iNk} = V \quad (5)$$

$$\sum_{i=0}^N \sum_{j=0}^N d_{ij} \cdot x_{ijk} \leq D, k \in U \quad (6)$$

$$z_{ik} - z_{jk} + D \cdot x_{ijk} \leq D - d_{ij}, i \neq j \in T, k \in U \quad (7)$$

$$0 \leq z_{ik} \leq D - d_{i0}, i \in \{1, 2, \dots, N\}, k \in U \quad (8)$$

Where  $D$  is the maximum flight distance of UAVs. The objective function (1) aims to maximize the total number of executed tasks. Constraints (2) to (4) ensure that each task can be executed once at most. Constraint (5) restricts all the UAVs must depart from and return to the depot. Constraint (6) is the constraint of flight distance for each UAV. Subtours are eliminated with constraints (7) and (8).

#### B. Divide and Conquer-based Framework

To reduce the complexity of the large-scale multi-UAV task scheduling problem, we divide it into two subproblems of task allocation and UAV route planning based on the divide and conquer principle. As shown in Fig. 1, the DCF which comprises an upper layer and a low layer is constructed. According to the initial scenario, several task subsets are generated and allocated to the corresponding UAVs in the upper layer. Afterward, each UAV performs the UAV route planning process in the lower layer based on its allocated tasks.

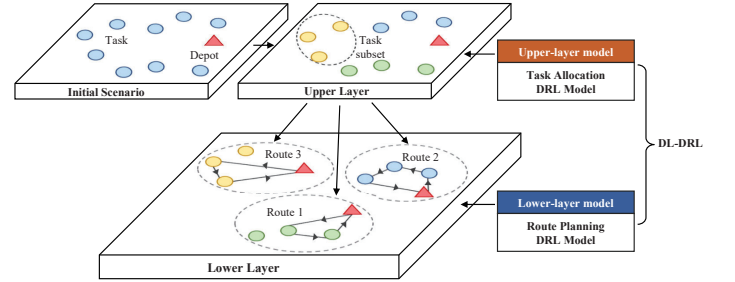


Fig. 1. The divide and conquer-based task scheduling framework

Based on DCF, we propose a DR-DRL approach for solving the multi-UAV task scheduling problem effectively, in which two different DRL models are used in the upper and lower layers, respectively. In the upper layer, we construct a task allocation DRL model called upper-layer model to select the appropriate UAV for each task. In the lower layer, we use a route planning DRL model called lower-layer model to plan the viable route for each UAV. Through the combination of these two layers, the multi-UAV task scheduling problem can be solved and this framework can be used to other COP, such as M-TSP [44], OP [45], etc.

### IV. METHODOLOGY

In this section, we introduce our DL-DRL approach for solving the multi-UAV task scheduling problem. The DCF decomposes the original problem into task allocation and UAV route planning subproblems. For the task allocation subproblem, we construct the Markov decision process (MDP) firstly and then design the policy networks of the task allocation DRL model. For the UAV route planning subproblem, we use

one of the state-of-the-art DRL models, AM [26], to solve this subproblem. Furthermore, we propose a novel training strategy that consists of pre-training, intensive training, and alternate training processes for effective model training.

### A. Markov Decision Process

The procedure of the task allocation can be deemed as a sequential decision-making process, where the allocation scheme is generated step by step. We model such a process as an MDP which includes state space  $S$ , action space  $A$ , state transition rule  $P$ , and reward  $R$ , i.e.,  $MDP = (S, A, P, R)$ . The detailed definition of our MDP is as follows.

**State:** the state includes the information of UAV and task, i.e.,  $s_t = (V_t, x_t) \in S$ , where  $V_t$  is the task set that have been allocated to UAVs at time step  $t$ , and  $x_t$  is the coordinates of the task that need to be allocated at time step  $t$ .  $V_t = \{U_t^1, U_t^2, \dots, U_t^v\} = \{\{u_1^1, u_2^1, \dots, u_t^1\}, \{u_1^2, u_2^2, \dots, u_t^2\}, \dots, \{u_1^v, u_2^v, \dots, u_t^v\}\}$ , where  $U_t^i$  and  $u_t^i$  is the task set and the corresponding task information of the UAV  $i$  at time step  $t$ .

**Action:** the action is to allocate the current task to a UAV. The action can be formulated as  $a_t = (u_i, x_t) \in A$ , i.e., the task  $x_t$  will be allocated to the UAV  $i$  at time step  $t$ .

**State transition rule:** the state transition rule will transmit the state  $s_t$  to  $s_{t+1}$  after implement the action  $a_t$ , i.e.,  $s_{t+1} = (V_{t+1}, x_{t+1}) = P(s_t, a_t)$ . The elements in  $V_t$  are updated as follows.

$$U_{t+1}^j = \begin{cases} [U_t^j, x_t^j], j = i \\ [U_t^j, u_t^j], \text{otherwise} \end{cases} \quad (9)$$

Where  $[,]$  is the concatenation of two vectors, and  $u_t^j$  is the last element of  $U_t^j$  (i.e., the allocated task for UAV  $j$  at time step  $t$ ). We use this update strategy to ensure the same dimension of each UAV's allocated task set, which is convenient for the batch training.

**Reward:** we define the reward as the total number of tasks executed by all UAVs to satisfies the objective function in the mathematical model, i.e.,  $R = \sum_{i=1}^v T_i$ , where  $T_i$  is the number of tasks executed by UAV  $i$ , which is obtained from the route planning DRL model.

### B. Architecture of Policy Networks

Compared to conventional Long Short-Term Memory and Gated Recurrent Unit models, the Transformer model has the stronger capability to handle sequential data and is extensively employed in domains of natural language processing [46], image recognition [47], recommendation system [48], and combination optimization [41]. Regarding the task allocation, we propose a Transformer-style policy network which is constructed by an encoder and a decoder to allocate tasks to the appropriate UAV. Fig. 2 shows the overall architecture of our policy network.

**Encoder:** The inputs of the encoder are the coordinates of tasks and the depot. We use a linear projection and attention layers to better extract the input features. Specifically, the raw inputs are first embedded to  $d_h$ -dimensional node embeddings

$h^0$  via the linear projection, where  $d_h = 128$ . The node embeddings  $h^0$  are updated to  $h^L$  through  $L$  attention layers. Each attention layer consists of a multi-head attention (MHA) sublayer and a feed-forward (FF) sublayer.

In the  $l$ -th attention layer, the input  $h^{l-1} = \{h_1^{l-1}, h_2^{l-1}, \dots, h_N^{l-1}\}$  is the output of the last attention layer, where  $N$  is the number of tasks. To engender the output of the  $l$ -th attention layer  $h^l$ , the MHA sublayer processes  $h^{l-1}$  via the multi-head self-attention mechanism first. For each head, the vectors of *query*, *key*, and *value* are obtained by multiplying  $h^{l-1}$  and the corresponding trainable parameters. Then, we calculate the attention value  $Y_{lm}$  according to Eq. (11), where  $m = \{1, 2, \dots, M\}$  and  $\dim = \frac{128}{M}$ . Here, the number of heads  $M = 8$ . After the calculation of attention values, we concatenate the attention values of all heads and project it to a  $d_h$ -dimensional vector, and then obtain the node embeddings  $\hat{h}^l$  through skip-connection and batch normalization operations. Afterward,  $\hat{h}^l$  is fed to the FF sublayer which is constructed by two linear projections and a ReLU activation function. Skip-connection and batch normalization are also used for the output of the FF sublayer to generate the output of the  $l$ -th attention layer. The process can be formulated as follows.

$$q_{lm} = W_{lm}^Q h^{l-1}, k_{lm} = W_{lm}^K h^{l-1}, v_{lm} = W_{lm}^V h^{l-1} \quad (10)$$

$$Y_{lm} = \text{softmax} \left( \frac{q_{lm}^T k_{lm}}{\sqrt{\dim}} \right) v_{lm} \quad (11)$$

$$\text{MHA}(h^{l-1}) = [Y_{11}; Y_{12}; \dots; Y_{1M}] W_l^O \quad (12)$$

$$\hat{h}^l = \text{BN}^l(h^{l-1} + \text{MHA}(h^{l-1})) \quad (13)$$

$$h^l = \text{BN}^l(\hat{h}^l + \text{FF}(\hat{h}^l)) \quad (14)$$

Where  $W_{lm}^Q, W_{lm}^K, W_{lm}^V \in \mathbb{R}^{M \times 128 \times \dim}$  and  $W_l^O \in \mathbb{R}^{128 \times 128}$  are trainable parameters in the  $l$ -th attention layer.

**Decoder:** The output of the encoder  $h^L$  is the input of the decoder. In the decoding process, we select an appropriate UAV for the current task at a timestep and add the task information (i.e., node embedding) to the UAV's task subset. Specifically, we construct an allocation feature context for each UAV.  $\tilde{C}_t^i = [h_0^{iL}, h_1^{iL}, \dots, h_j^{iL}, \dots, h_{t-1}^{iL}]$  is the allocation feature context of UAV  $i$ , where  $h_j^{iL}$  is the node embedding of the task that is allocated to UAV  $i$  at the timestep  $j$ . We aggregate all the allocation feature contexts via a max pooling and concatenate them to obtain a task subset context  $\tilde{C}_t^{TS}$ . After that, a linear projection and a 512-dimensional FF layer are used to generate the task subset embeddings  $H_t^{TS}$ . Finally, we concatenate  $H_t^{TS}$  and node embedding of the current task  $h_t^L$ , and then feed it to a linear projection and the softmax function to obtain the output probability. Based on the probability, we select the UAV using one of two decoding rules: greedy or sampling. Greedy rule always selects the UAV with the highest probability, while the sampling rule selects the UAV according to its probability. We show these steps as follows.



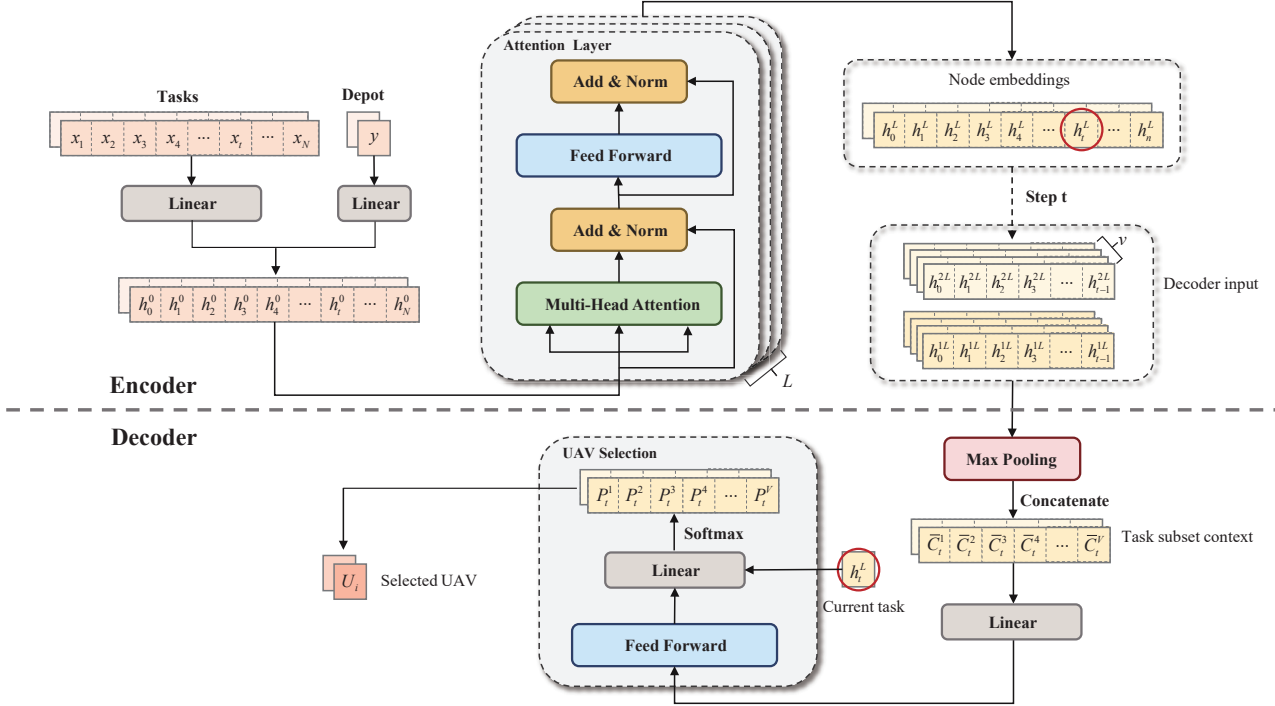


Fig. 2. Architecture of policy network for multi-UAV task allocation

$$\bar{C}_t^i = \max(\hat{C}_t^i) \quad (15)$$

$$\hat{C}_t^{TS} = [\bar{C}_t^1, \bar{C}_t^2, \dots, \bar{C}_t^V] \quad (16)$$

$$H_t^{TS} = FF(W_1 \hat{C}_t^{TS} + b_1) \quad (17)$$

$$p_t = \text{softmax}(W_2 [H_t^{TS}, h_t^N] + b_2) \quad (18)$$

Where,  $W_1$ ,  $b_1$ ,  $W_2$  and  $b_2$  are trainable parameters.

### C. Interactive Training Strategy

Based on DCF, we use two different DRL models (i.e., upper-layer model and lower-layer model) to solve the task allocation and UAV route planning subproblems, respectively. These two models are interdependent as the input data of the lower-layer model is generated by the upper-layer model, and the reward of the upper-layer model is obtained via the lower-layer model. Therefore, we design an ITS that consists of pre-training, intensive training, and alternate training processes for effective model training. The pseudocode of our ITS is presented in Algorithm 1. As for the training algorithm, we use the rollout-based REINFORCE which performs well in training Transformer-style DRL models [41], [42], [49]. The gradient of the loss function is calculated below:

$$\nabla \mathcal{L}(\theta|s) = \mathbb{E}_{p_\theta(\pi|s)}[(L(\pi) - L(\pi^{BL}))\nabla_\theta \log p_\theta(\pi|s)] \quad (19)$$

Where  $\theta$  is the parameters of policy network  $\pi$ ,  $L(\pi)$  and  $L(\pi^{BL})$  are the cost function (i.e., the negative number of executed tasks) of the training model and the baseline model, respectively. We use the greedy decoding strategy in

the baseline model. Its parameters will be replaced by that of the current training model if the current model has significant improvement according to the paired t-test. In this way, the loss variance will effectively decrease during the whole training process.

---

#### Algorithm 1 Interactive training strategy

---

**Input:** the number of pre-training epochs  $E_p$ ; the number of models training epochs  $E$ ; the number of intensive training epochs  $E_t$ ; the number of continuous training epochs  $E_c$ ;

**Initialization:** the training datasets of the pre-training, intensive training and alternate training processes; the parameters of networks in the upper-layer model  $M_u$  and the lower-layer model  $M_l$ ;

```

1: for  $i = 1, 2, \dots, E_p$  do
2:   REINFORCE( $M_l$ )
3: end for
4: for  $epoch = 1, 2, \dots, E$  do
5:   REINFORCE( $M_u$ )
6:   if  $epoch < E_t$  then
7:     if  $epoch \% E_c == 0$  then
8:       Generate and shuffle the training datasets of the
9:       lower-layer model  $M_l$ 
10:      REINFORCE( $M_l$ )
11:    end if
12:  else
13:    Generate and shuffle the training datasets of the
14:    lower-layer model  $M_l$ 
15:    REINFORCE( $M_l$ )
16:  end if
17: end for

```

---

At the beginning of the models training, we construct the pre-training process for the lower-layer model (lines 1-3), where the model is trained for  $E_p$  epochs. This process can make the lower-layer model learn an initial route planning policy to support the following training process. After the pre-training process, the higher-layer and lower-layer models are trained interactively, with the lower-layer model fixed when training the higher-layer model and vice versa. During the interactive training, there are two processes: 1) intensive training process (lines 5-10), in which the higher-layer model is trained for  $E_c$  epochs continuously and the lower-layer model is trained for the same number of epochs; 2) alternate training process (line 5 and lines 12-13), in which the higher-layer model and lower-layer model are trained for one epoch iteratively. In addition, after generating training datasets for the lower-layer model, we should shuffle them to ensure that the input data of the lower-layer model is independent and identically distributed.

## V. EXPERIMENT

To evaluate the performance of our DL-DRL for solving the multi-UAV task scheduling problem, we conduct extensive experiments to compare it with several state-of-the-art baselines. The effectiveness of the ITS is also evaluated by comparing the convergence performance of the model with different training strategies. Besides, we further test the DL-DRL on instances with larger-scale tasks to verify its generalization performance.

### A. Experiment Settings

Following the data generation in [26], [49], [50], we randomly sample the location of tasks and depot in the uniform distribution  $U[0, 1]$ . Two scenarios with four and six UAVs are considered (named U4 and U6, respectively), and each scenario is classified into small-scale situations, medium-scale situations, and large-scale situations according to the problem size. Instances with 80 and 100 tasks, 150 and 200 tasks, and 300 and 500 tasks are generated in small-scale, medium-scale, and large-scale situations, respectively. In this way, we can evaluate the overall performance of the approaches across a range of task scales and UAV quantities. Furthermore, the maximum flight distance of the UAV is set to 2.0 in all instances, and the number of attention layers in the encoder  $L = 3$ .

In our DL-DRL, both the upper-layer and lower-layer DRL models are trained for 100 epochs (except for pre-training) with 1,280,000 randomly generated instances per epoch. The batch size is set to 512 for the model training in small-scale and medium-scale situations, while 256 for the large-scale situation due to the memory constraint. We utilize the Adam optimizer to update the network parameters in DRL models with an initial learning rate of  $10^{-4}$ . In the upper-layer model, the norm of grads is clipped within 3.0 and the decaying parameter of the learning rate is set to 0.995, while the same hyperparameters are set to 1.0 and 1.0 in the lower-layer model. In ITS, the lower-layer model is pre-trained for five epochs (i.e.,  $E_p = 5$ ), and different parameter settings are used for U4 and U6 scenarios during the interactive training.

For the U4 scenario, The number of intensive training epochs  $E_t$  is set to 8, and the number of continuous training epochs  $E_c$  is set to 4, whereas  $E_t$  and  $E_c$  are set to 12 and 6 for the U6 scenario, respectively. In this way, we avoid the violent shaking of models which results from the imbalanced amount of training data of the upper-layer and lower-layer models throughout the training process. Our code is publicly available <sup>1</sup>.

### B. Learning Performance

Following the experiment settings, DL-DRL models are trained for U4 and U6 scenarios in different situations. For the model training in the small-scale and medium-scale situations, a single RTX 3090 GPU with 24GB memory is used, whereas two GPUs are used for the model training in the large-scale situation. The time consuming of the upper-layer model, lower-layer model, and data generation for the lower-layer model per epoch is presented in Table I, where “U4-80” means the model is trained for the U4 scenario with 80 tasks. After the model training process, we visualize the results of trained models in Fig. 3, where all the problem instances are randomly generated on seed 1234. From Fig. 3, it is obvious that the trained models can obtain reasonable solutions in their corresponding situations.

TABLE I  
TIME CONSUMING OF MODEL TRAINING FOR EACH EPOCH

| Situation | Upper-layer | Lower-layer | Data generation |
|-----------|-------------|-------------|-----------------|
| U4-80     | 20.55min    | 6.12 min    | 13.02 min       |
| U4-100    | 24.58 min   | 6.58 min    | 14.06 min       |
| U4-150    | 32.15 min   | 8.78 min    | 19.30 min       |
| U4-200    | 55.98 min   | 11.15 min   | 22.77 min       |
| U4-300    | 123.70min   | 23.32min    | 45.50min        |
| U4-500    | 206.37min   | 31.28min    | 84.95min        |
| U6-80     | 25.92min    | 6.03min     | 17.19min        |
| U6-100    | 30.65min    | 6.63min     | 18.58min        |
| U6-150    | 49.25min    | 8.55 min    | 23.35 min       |
| U6-200    | 72.10 min   | 11.32 min   | 29.15 min       |
| U6-300    | 155.55min   | 22.92min    | 58.49min        |
| U6-500    | 262.43min   | 31.27min    | 116.02min       |

Furthermore, we plot the learning curves of the upper-layer and lower-layer models during the training process in Fig. 4. The vertical coordinate is the cost value, which is the negative of the total number of executed tasks. Hence, the lower value means better performance. From Fig. 4, we can observe two phenomena: 1) For the lower-layer model, it converges fast during the pre-training process, while the cost value dramatically increases at the beginning of the interactive training. This phenomenon is caused by the difference in the training data distribution between pre-training and interactive training, i.e., the former is a uniform distribution and the latter is prior unknown and generated by the upper-layer model; 2) For the upper-layer model, the cost value has several rapid dropping procedures in the early training period. As the lower-layer model learns a better policy based on the training data generated by the upper-layer model during the intensive training process, the upper-layer model quickly

<sup>1</sup>[https://faculty.csu.edu.cn/guohuawu/zh\\_CN/zdylm/193832/list/index.htm](https://faculty.csu.edu.cn/guohuawu/zh_CN/zdylm/193832/list/index.htm)

obtains good performance. In addition, the learning curves have more obvious fluctuation in the large-scale situation due to the expansion of the state and action space. The DL-DRL models finally converge stably, demonstrating our models have learned a valid policy.

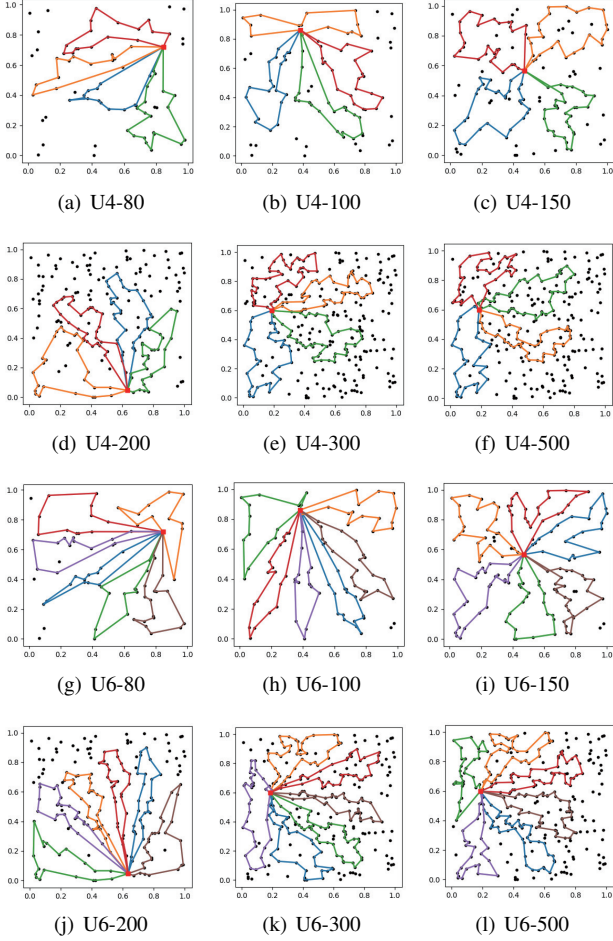


Fig. 3. The solution visualization of trained models in corresponding situations

### C. Comparison Analysis

We compare the proposed DL-DRL with the exact, heuristic, and learning-based baselines as follows.

- 1) Gurobi<sup>2</sup>: A commercial mathematical optimizer;
- 2) OR-Tools<sup>3</sup>: A state-of-the-art heuristic optimization solver developed by Google;
- 3) K-means+VND: Based on DCF, using K-means and variable neighborhood descent (VND)<sup>4</sup> for task allocation and UAV route planning, respectively;
- 4) K-means+AM: Based on DCF, using K-means and AM [26] for task allocation and UAV route planning, respectively;
- 5) DL-DRL-I: An approach with the same network architecture as our DL-DRL, while it trains the DRL models independently rather than using the ITS.

As an exact optimizer, Gurobi is time-consuming for solving multi-UAV task scheduling problems even with small problem sizes. Therefore, for each problem size, we evaluate Gurobi with 30 testing instances and limit the maximum solving time to 1800s, whereas other baselines are tested with 500 testing instances without time limitation. All of the baselines are implemented in Python. The Intel Xeon Gold 5218R with 20 cores and RTX 3090 GPU are utilized for traditional and learning-based approaches, respectively.

Regarding the upper-layer model of DL-DRL and DL-DRL-I, both greedy and sampling strategies are used in the decoder during the testing process. According to Eq. (18), the greedy strategy always selects the UAV with the highest probability, whereas the sampling strategy engenders 128 solutions by sampling and reports the best. The experimental results in U4 and U6 scenarios are presented in Table II and Table III, respectively. These tables give the average objective value (Obj.), gap, and the average computation time of all methods for each situation. Here, the Obj. is the total number of executed tasks on average, and the gap is defined as the normalized distance between an average objective value  $Obj$  and the best one  $Obj_{best}$  among all methods.

$$Gap = \frac{Obj_{best} - Obj}{Obj_{best}} \times 100\% \quad (20)$$

We can see from Tables II and III that Gurobi is unable to obtain the optimal solution, even for instances with only 80 tasks. At the same time, our DL-DRL always obtains better solutions than Gurobi no matter which decoding strategy is used.

For the experimental results in the U4 scenario from Table II, it is obvious that our DL-DRL (Greedy) outperforms K-means+VND and K-means+AM in terms of Obj. and computation time, with the advantage becoming more significant as the task scale increases. Regarding the DL-DRL and DL-DRL-I, which have the same network architectures, the sampling strategy always generates better solutions than the greedy strategy with only a slight increase in computation time. Our DL-DRL (Greedy) outperforms DL-DRL-I (Greedy) and the DL-DRL-I (Sampling), which indicates the effectiveness of the ITS in improving the model's performance. Compared to the state-of-the-art heuristic (i.e., OR-Tools), our DL-DRL (Sampling) is slightly worse in terms of Obj. in the problems with the fewest tasks. However, with the increasing task scale, DL-DRL (Sampling) outperforms OR-Tools in terms of both Obj. and computation time, which demonstrates the strong capability of our approach in addressing large-scale task scheduling problems of multi-UAV.

In Table III, a similar observation can also be found that our DL-DRL (Greedy and Sampling) outperforms K-means+VND, K-means+AM, DL-DRL-I (Greedy), and DL-DRL-I (Sampling). Meanwhile, our DL-DRL (Sampling) has a competitive performance against OR-Tools. In situations with no more than 150 tasks, DL-DRL (Sampling) is slightly inferior to OR-Tools in terms of Obj., but its computation time is much shorter. When the number of tasks increases, our DL-DRL (Sampling) outperforms OR-Tools in terms of Obj. and

<sup>2</sup><https://www.gurobi.com/>

<sup>3</sup><https://developers.google.cn/optimization/>

<sup>4</sup>[https://github.com/nchlpz/VND\\_TOP](https://github.com/nchlpz/VND_TOP)

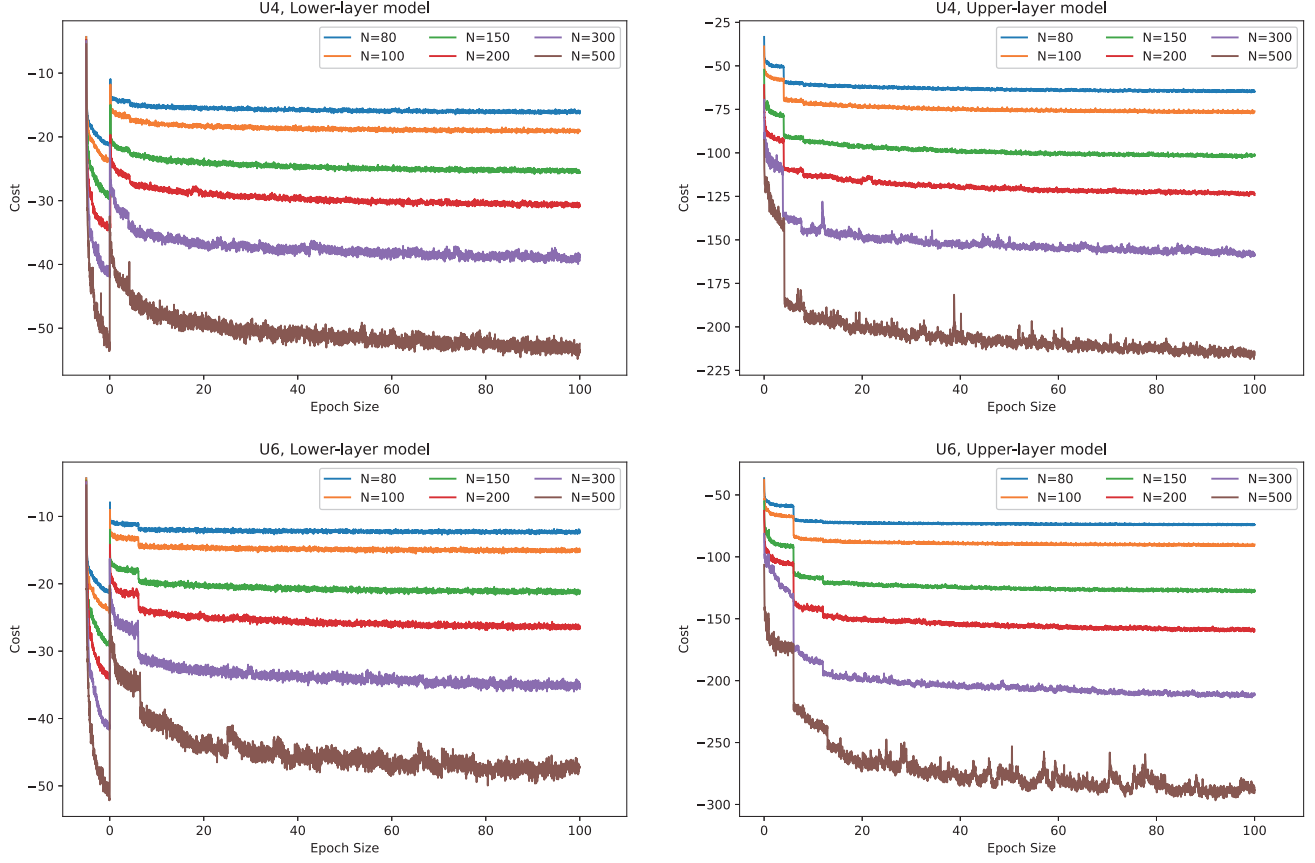


Fig. 4. The convergence curves of DL-DRL

TABLE II  
THE COMPARATIVE EXPERIMENT RESULTS IN THE U4 SCENARIO

| Method              | N=80         |             |        | N=100        |             |        | N=150         |             |        | N=200         |             |        | N=300         |             |         | N=500         |             |         |
|---------------------|--------------|-------------|--------|--------------|-------------|--------|---------------|-------------|--------|---------------|-------------|--------|---------------|-------------|---------|---------------|-------------|---------|
|                     | Obj.         | Gap         | Time   | Obj.         | Gap         | Time   | Obj.          | Gap         | Time   | Obj.          | Gap         | Time   | Obj.          | Gap         | Time    | Obj.          | Gap         | Time    |
| Gurobi(1800s)       | 56.77        | 11.72       | 1800s  | 65.37        | 14.96       | 1800s  | 75.10         | 27.53       | 1800s  | 74.37         | 40.54       | 1800s  | 61.77         | 61.21       | 1800s   | 69.33         | 68.83       | 1800s   |
| DL-DRL (Greedy)     | 63.67        | 0.98        | 0.03s  | 76.20        | 0.87        | 0.03s  | 102.57        | 1.03        | 0.03s  | 124.07        | 0.80        | 0.04s  | 158.50        | 0.46        | 0.04s   | 218.93        | 1.57        | 0.05s   |
| DL-DRL (Sampling)   | <b>64.30</b> | <b>0.00</b> | 0.11s  | <b>76.87</b> | <b>0.00</b> | 0.13s  | <b>103.63</b> | <b>0.00</b> | 0.21s  | <b>125.07</b> | <b>0.00</b> | 0.29s  | <b>159.23</b> | <b>0.00</b> | 0.44s   | <b>222.43</b> | <b>0.00</b> | 0.89s   |
| OR-tools            | <b>65.44</b> | <b>0.00</b> | 0.38s  | 76.83        | 0.72        | 0.60s  | 101.34        | 2.01        | 1.25s  | 121.21        | 2.87        | 2.00s  | 155.60        | 2.84        | 4.41s   | 209.00        | 5.19        | 11.71s  |
| K-means+VND         | 56.45        | 13.74       | 39.20s | 62.65        | 19.05       | 53.31s | 74.11         | 28.34       | 67.79s | 81.47         | 34.72       | 96.51s | 94.86         | 40.77       | 164.32s | 113.88        | 48.34       | 316.45s |
| K-means+AM          | 59.68        | 8.79        | 0.12s  | 69.35        | 10.38       | 0.13s  | 87.14         | 15.74       | 0.39s  | 100.85        | 19.19       | 0.38s  | 124.77        | 22.09       | 0.54s   | 154.71        | 29.82       | 0.66s   |
| DL-DRL-I (Greedy)   | 58.78        | 10.18       | 0.00s  | 67.61        | 12.63       | 0.00s  | 90.28         | 12.71       | 0.00s  | 106.68        | 14.51       | 0.00s  | 134.42        | 16.06       | 0.00s   | 170.20        | 22.79       | 0.01s   |
| DL-DRL-I (Sampling) | 61.32        | 6.29        | 0.12s  | 71.14        | 8.07        | 0.15s  | 96.41         | 6.77        | 0.21s  | 115.34        | 7.57        | 0.30s  | 145.41        | 9.20        | 0.42s   | 188.78        | 14.36       | 0.84s   |
| DL-DRL (Greedy)     | 64.67        | 1.17        | 0.00s  | 76.67        | 1.01        | 0.00s  | 102.25        | 1.12        | 0.00s  | 123.76        | 0.83        | 0.00s  | 159.25        | 0.56        | 0.00s   | 216.66        | 1.72        | 0.01s   |
| DL-DRL (Sampling)   | 65.28        | 0.24        | 0.11s  | <b>77.39</b> | <b>0.00</b> | 0.14s  | <b>103.42</b> | <b>0.00</b> | 0.21s  | <b>124.79</b> | <b>0.00</b> | 0.27s  | <b>160.14</b> | <b>0.00</b> | 0.41s   | <b>220.44</b> | <b>0.00</b> | 0.85s   |

computation time. Considering the results of both tables, DL-DRL has a better performance than Gurobi, K-means+VND, K-means+AM, and DL-DRL-I. It also outperforms OR-Tools in large-scale situations and performs competitively well against OR-Tools with a shorter computation time when the problem size is small.

#### D. Ablation Study on ITS

To investigate the influence of different processes in ITS, we construct the ablation study in the scenario with 4 UAVs and 80 tasks. Different training strategies are used for the model training, we record the cost value during the training process. Fig. 5 shows the cost value curves throughout the training process, where “ITS/X” denotes the training strategy that removes the “X” process from our ITS. For example, “ITS/intensive” denotes the training strategy that only includes the pre-training and alternate training processes.

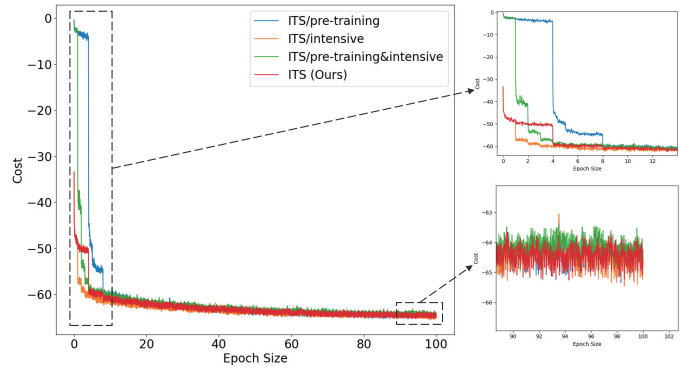


Fig. 5. Ablation study results of ITS

From Fig. 5, we can observe that the pre-trained models perform better at the beginning of the interactive training (e.g., “ITS” versus “ITS/pre-training”). The intensive training process slows convergence speed by several epochs (e.g.,

TABLE III  
THE COMPARATIVE EXPERIMENT RESULTS IN THE U6 SCENARIO

| Method              | Obj.         | N=80<br>Gap | Time   | Obj.         | N=100<br>Gap | Time   | Obj.          | N=150<br>Gap | Time   | Obj.          | N=200<br>Gap | Time    | Obj.          | N=300<br>Gap | Time    | Obj.          | N=500<br>Gap | Time    |
|---------------------|--------------|-------------|--------|--------------|--------------|--------|---------------|--------------|--------|---------------|--------------|---------|---------------|--------------|---------|---------------|--------------|---------|
| Gurobi(1800s)       | 64.37        | 12.47       | 1800s  | 73.23        | 19.73        | 1800s  | 79.63         | 38.56        | 1800s  | 73.27         | 54.20        | 1800s   | 79.80         | 62.76        | 1800s   | 61.00         | 79.65        | 1800s   |
| DL-DRL (Greedy)     | 73.10        | 0.59        | 0.03s  | 90.27        | 1.06         | 0.03s  | 128.23        | 1.06         | 0.03s  | 157.37        | 1.63         | 0.04s   | 211.93        | 1.09         | 0.05s   | 295.40        | 1.45         | 0.06s   |
| DL-DRL (Sampling)   | <b>73.53</b> | <b>0.00</b> | 0.15s  | <b>91.23</b> | <b>0.00</b>  | 0.19s  | <b>129.60</b> | <b>0.00</b>  | 0.27s  | <b>159.97</b> | <b>0.00</b>  | 0.37s   | <b>214.27</b> | <b>0.00</b>  | 0.58s   | <b>299.73</b> | <b>0.00</b>  | 1.18s   |
| OR-tools            | <b>75.42</b> | <b>0.00</b> | 0.50s  | <b>92.42</b> | <b>0.00</b>  | 0.85s  | <b>130.09</b> | <b>0.00</b>  | 2.21s  | 160.79        | 0.03         | 3.73s   | 212.46        | 1.40         | 8.00s   | 292.10        | 1.84         | 20.84s  |
| K-means+VND         | 68.59        | 9.05        | 20.65s | 81.38        | 11.95        | 36.32s | 104.49        | 19.67        | 88.55s | 117.54        | 26.92        | 151.36s | 135.89        | 36.94        | 221.86s | 162.25        | 45.48        | 382.92s |
| K-means+AM          | 68.31        | 9.43        | 0.14s  | 82.12        | 11.14        | 0.16s  | 110.14        | 15.33        | 0.19s  | 136.05        | 15.41        | 0.23s   | 165.10        | 23.38        | 0.28s   | 193.21        | 35.07        | 0.32s   |
| DL-DRL-I (Greedy)   | 68.91        | 8.63        | 0.00s  | 82.60        | 10.63        | 0.00s  | 110.53        | 15.04        | 0.00s  | 137.09        | 14.76        | 0.00s   | 178.61        | 17.11        | 0.01s   | 201.76        | 32.20        | 0.01s   |
| DL-DRL-I (Sampling) | 71.01        | 5.85        | 0.13s  | 85.43        | 7.56         | 0.19s  | 116.29        | 10.60        | 0.26s  | 144.78        | 9.98         | 0.35s   | 191.23        | 11.25        | 0.57s   | 226.21        | 23.99        | 1.07s   |
| DL-DRL (Greedy)     | 74.00        | 1.88        | 0.00s  | 90.72        | 1.84         | 0.00s  | 128.18        | 1.46         | 0.01s  | 158.96        | 1.16         | 0.01s   | 213.27        | 1.02         | 0.01s   | 293.18        | 1.48         | 0.01s   |
| DL-DRL (Sampling)   | 74.32        | 1.46        | 0.14s  | 91.53        | 0.96         | 0.17s  | 129.59        | 0.38         | 0.24s  | <b>160.83</b> | <b>0.00</b>  | 0.35s   | <b>215.47</b> | <b>0.00</b>  | 0.53s   | <b>297.59</b> | <b>0.00</b>  | 1.05s   |

“ITS/pre-training” versus “ITS/pre-training&intensive”). However, for each epoch, the time consuming of the training data generation for the lower-layer model is much longer than that of the lower-layer model training, as listed in Table I. The intensive training process can skip the time-consuming procedure of training data generation for the lower-layer model, which efficiently reduces the whole training time. With respect to the final performance of models, ITS (ours) has a similar performance to ITS/intensive, whereas the whole training time of ITS (ours) is much shorter than that of ITS/intensive due to the skip of lower-layer training data generation. From Table I, the training data generation for the lower-layer model can consume several hours in large-scale situations. Thus, we use the ITS which has an intensive training process to accelerate the entire training process with no performance decrease.

#### E. Generalization on Larger-scale Problems

To verify the generality of our DL-DRL, we use the trained models to solve larger-scale problems and compare them to several baselines. The DL-DRL (Sampling) is used to evaluate the generalization performance since the sampling strategy can obtain a better solution with only a slight increase in computation time. We first generalize the trained DL-DRL models from subsection V-B to solve the problem with a larger task scale. The results are shown in Fig. 6 and 7, where the horizontal coordinate refers to the problems that need to be solved, and the label denotes a certain trained model. For example, N-100 represents the problem with 100 tasks, and U4-80 represents the DL-DRL model trained for the U4 scenario with 80 tasks.

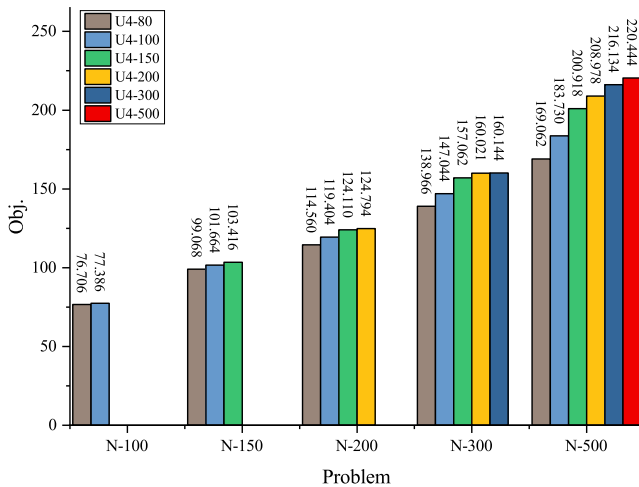


Fig. 6. Generalization performance of trained model in U4

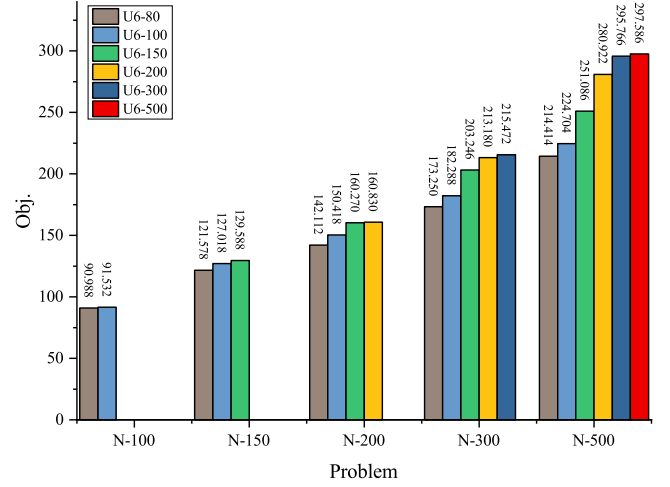


Fig. 7. Generalization performance of trained model in U6

From Fig. 6, we can see that the model trained for a certain problem size performs best in the corresponding problem compared to those trained for other problem sizes. However, they still outperform K-means+VND, K-means+AM, and DL-DRL-I except for U4-80 in solving problems with task quantity greater than or equal to 200 and U4-100 in solving the problem with 500 tasks. In addition, we notice that the model trained for proximal problem sizes performs better than that of different ones (e.g., U4-200 and U4-300 perform better than U4-80, U4-100, and U4-150 in solving problems with 500 tasks). This phenomenon might be caused by the difference in data distribution, as the proximal problem sizes may lead to similar data distributions of task location. A similar observation can also be found in Fig. 7, where the models trained for other problem sizes outperform K-means+VND and K-means+AM and are competitive against DL-DRL-I.

For multi-UAV task scheduling problems with larger-scale tasks such as 1000, it is intractable to train the model from scratch, making the generalizability of trained models for larger-scale problems significant. To further investigate the generalization performance of the proposed DL-DRL, we construct the experiments in the scenarios with 600, 700, 800, 900, and 1000 tasks. Comparing the DL-DRL model trained for 500 tasks with baselines, i.e., OR-Tools, K-means+AM, and DL-DRL-I. Table IV shows the testing results of our DL-DRL and baselines, including Obj. and average computation time. Compared to K-means+AM and DL-DRL-I, the significant superiority of DL-DRL can be observed in terms of Obj. with a slightly longer computation time. Regarding OR-Tools, the computation time is rapidly increasing as the task scale



TABLE IV  
GENERALIZATION RESULTS OF DL-DRL AND BASELINES

|    | Method     | N=600         |              | N=700         |              | N=800         |              | N=900         |              | N=1000        |              |
|----|------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|
|    |            | Obj.          | Time         | Obj.          | Time         | Obj.          | Time         | Obj.          | Time         | Obj.          | Time         |
| U4 | OR-tools   | 231.79        | 16.93s       | 251.65        | 22.15s       | 271.20        | 28.50s       | 288.36        | 36.17s       | 305.64        | 43.01s       |
|    | K-means+AM | 165.73        | <b>0.27s</b> | 174.75        | <b>0.28s</b> | 181.54        | <b>0.29s</b> | 187.72        | <b>0.30s</b> | 192.43        | <b>0.31s</b> |
|    | DL-DRL-I   | 206.46        | 1.06s        | 220.48        | 1.31s        | 232.94        | 1.61s        | 241.13        | 1.93s        | 248.81        | 2.28s        |
|    | DL-DRL     | <b>245.67</b> | 1.07s        | <b>267.88</b> | 1.33s        | <b>287.36</b> | 1.62s        | <b>303.83</b> | 1.95s        | <b>321.58</b> | 2.34s        |
| U6 | OR-tools   | 326.37        | 30.45s       | 356.94        | 39.70s       | 386.45        | 52.30s       | 411.60        | 65.73s       | 437.17        | 79.83s       |
|    | K-means+AM | 209.12        | <b>0.35s</b> | 221.79        | <b>0.37s</b> | 230.78        | <b>0.38s</b> | 238.52        | 0.39s        | 244.79        | <b>0.39s</b> |
|    | DL-DRL-I   | 254.55        | 1.47s        | 277.81        | 1.78s        | 295.52        | 2.21s        | 305.69        | 2.67s        | 316.94        | 3.24s        |
|    | DL-DRL     | <b>334.28</b> | 1.43s        | <b>367.61</b> | 1.84s        | <b>394.90</b> | 2.28s        | <b>419.26</b> | 2.74s        | <b>443.20</b> | 3.32s        |

grows, and our DL-DRL outperforms it in terms of both Obj. and computation time. In all the scenarios, our DL-DRL achieves the best overall performance among all methods, which demonstrates its great generalization performance.

## VI. CONCLUSION

In this paper, we investigate the DRL approach for the large-scale multi-UAV task scheduling problem. Based on DCF, the original problem is divided into task allocation and UAV route planning subproblems, and we proposed the DL-DRL approach. In our DL-DRL, two different DRL models are constructed and used for solving the two subproblems, respectively. Considering the intrinsic relationship of subproblems, the ITS is designed for model training, which comprises pre-training, intensive training, and alternate training processes. The experimental results show that our DL-DRL has the overall best performance compared to exact, heuristic, and learning-based baselines, especially in large-scale situations. Furthermore, the generality of the DL-DRL performs pretty well in larger-scale problems and the effectiveness of our ITS is also verified by the ablation study. With respect to future studies, we notice that the model is trained and tested in a uniform distribution, and we will improve our approach to face situations with different data distributions.

## VII. ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China under Grants 62073341 and the Fundamental Research Funds for the Central Universities of Central South University (No. 2022ZZTS0659).

## REFERENCES

- [1] Y. Liu, "An optimization-driven dynamic vehicle routing algorithm for on-demand meal delivery using drones," *Computers & Operations Research*, vol. 111, pp. 1–20, Nov. 2019.
- [2] D. Machovec, J. A. Crowder, H. J. Siegel, S. Pasricha, and A. A. Maciejewski, "Dynamic Heuristics for Surveillance Mission Scheduling with Unmanned Aerial Vehicles in Heterogeneous Environments," in *Advances in Artificial Intelligence and Applied Cognitive Computing*. Cham: Springer International Publishing, 2021, pp. 583–605.
- [3] A. Altan and R. Hacıoğlu, "Model predictive control of three-axis gimbal system mounted on UAV for real-time target tracking under external disturbances," *Mechanical Systems and Signal Processing*, vol. 138, p. 106548, Apr. 2020.
- [4] Z. Wang, L. Liu, T. Long, and Y. Wen, "Multi-UAV reconnaissance task allocation for heterogeneous targets using an opposition-based genetic algorithm with double-chromosome encoding," *Chinese Journal of Aeronautics*, vol. 31, no. 2, pp. 339–350, Feb. 2018.
- [5] G. Laporte and Y. Nobert, "A Cutting Planes Algorithm for the m-Salesmen Problem," *Journal of the Operational Research Society*, vol. 31, no. 11, pp. 1017–1023, Nov. 1980.
- [6] B. P. Gerkey and M. J. Mataric, "A Formal Analysis and Taxonomy of Task Allocation in Multi-Robot Systems," *The International Journal of Robotics Research*, vol. 23, no. 9, pp. 939–954, Sep. 2004.
- [7] A. Richards, J. Bellingham, M. Tillerson, and J. How, "Coordination and Control of Multiple UAVs," in *AIAA Guidance, Navigation, and Control Conference and Exhibit*. American Institute of Aeronautics and Astronautics, Aug. 2002.
- [8] M. Alighanbari and J. How, "Cooperative task assignment of unmanned aerial vehicles in adversarial environments," in *Proceedings of the 2005, American Control Conference, 2005*. Portland, OR, USA: IEEE, 2005, pp. 4661–4666.
- [9] E. J. Forsmo, E. I. Grøtli, T. I. Fossen, and T. A. Johansen, "Optimal search mission with Unmanned Aerial Vehicles using Mixed Integer Linear Programming," in *2013 International Conference on Unmanned Aircraft Systems (ICUAS)*, May 2013, pp. 253–259.
- [10] E. Edison and T. Shima, "Integrated task assignment and path optimization for cooperating uninhabited aerial vehicles using genetic algorithms," *Computers & Operations Research*, vol. 38, no. 1, pp. 340–356, Jan. 2011.
- [11] S. Gao, J. Wu, and J. Ai, "Multi-UAV reconnaissance task allocation for heterogeneous targets using grouping ant colony optimization algorithm," *Soft Computing*, vol. 25, no. 10, pp. 7155–7167, May 2021.
- [12] N. Geng, Z. Chen, Q. A. Nguyen, and D. Gong, "Particle swarm optimization algorithm for the optimization of rescue task allocation with uncertain time constraints," *Complex & Intelligent Systems*, vol. 7, no. 2, pp. 873–890, Apr. 2021.
- [13] L. Huo, J. Zhu, G. Wu, and Z. Li, "A Novel Simulated Annealing Based Strategy for Balanced UAV Task Assignment and Path Planning," *Sensors*, vol. 20, no. 17, p. 4769, Jan. 2020.
- [14] H. Liu, X. Li, G. Wu, M. Fan, R. Wang, L. Gao, and W. Pedrycz, "An iterative two-phase optimization method based on divide and conquer framework for integrated scheduling of multiple UAVs," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 9, pp. 5926–5938, Sep. 2021.
- [15] F. Zitouni, S. Harous, and R. Maamri, "A Distributed Solution to the Multi-robot Task Allocation Problem Using Ant Colony Optimization and Bat Algorithm," in *Advances in Machine Learning and Computational Intelligence*. Singapore: Springer, 2021, pp. 477–490.
- [16] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [17] Y. Zhu and D. Zhao, "Online Minimax Q Network Learning for Two-Player Zero-Sum Markov Games," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 3, pp. 1228–1241, Mar. 2022.
- [18] W. Zhao, J. P. Queralta, and T. Westerlund, "Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: A Survey," in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, 2020, pp. 737–744.
- [19] J. Luketina, N. Nardelli, G. Farquhar, J. Foerster, J. Andreas, E. Grefenstette, S. Whiteson, and T. Rocktäschel, "A Survey of Reinforcement Learning Informed by Natural Language," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*. Macao, China: International Joint Conferences on Artificial Intelligence Organization, Aug. 2019, pp. 6309–6317.
- [20] H. Lu, X. Zhang, and S. Yang, "A Learning-based Iterative Method for Solving Vehicle Routing Problems," in *International Conference on Learning Representations*, Sep. 2019.

- [21] J. Zheng, K. He, J. Zhou, Y. Jin, and C.-M. Li, "Combining Reinforcement Learning with Lin-Kernighan-Helsgaun Algorithm for the Traveling Salesman Problem," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 14, pp. 12 445–12 452, May 2021.
- [22] R. Gama and H. L. Fernandes, "A reinforcement learning approach to the orienteering problem with time windows," *Computers & Operations Research*, vol. 133, p. 105357, Sep. 2021.
- [23] Y. Ma, X. Hao, J. Hao, J. Lu, X. Liu, T. Xialiang, M. Yuan, Z. Li, J. Tang, and Z. Meng, "A Hierarchical Reinforcement Learning Based Optimization Framework for Large-scale Dynamic Pickup and Delivery Problems," in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 23 609–23 620.
- [24] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014.
- [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., Dec. 2017, pp. 6000–6010.
- [26] W. Kool, H. van Hoof, and M. Welling, "ATTENTION, LEARN TO SOLVE ROUTING PROBLEMS!" in *International Conference on Learning Representations*, 2019.
- [27] P.-L. Bacon, J. Harb, and D. Precup, "The Option-Critic Architecture," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, Feb. 2017.
- [28] O. Nachum, S. S. Gu, H. Lee, and S. Levine, "Data-Efficient Hierarchical Reinforcement Learning," in *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc., 2018.
- [29] C. Ramirez-Atencia, G. Bello-Organ, M. D. R.-Moreno, and D. Camacho, "Performance Evaluation of Multi-UAV Cooperative Mission Planning Models," in *Computational Collective Intelligence*. Cham: Springer International Publishing, 2015, pp. 203–212.
- [30] D. Mitchell, M. Corah, N. Chakraborty, K. Sycara, and N. Michael, "Multi-robot long-term persistent coverage with fuel constrained robots," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 1093–1099.
- [31] F. Mufalli, R. Batta, and R. Nagi, "Simultaneous sensor selection and routing of unmanned aerial vehicles for complex mission plans," *Computers & Operations Research*, vol. 39, no. 11, pp. 2787–2799, Nov. 2012.
- [32] F. Ye, J. Chen, Y. Tian, and T. Jiang, "Cooperative Task Assignment of a Heterogeneous Multi-UAV System Using an Adaptive Genetic Algorithm," *Electronics*, vol. 9, no. 4, p. 687, Apr. 2020.
- [33] K. Shang, S. Karungaru, Z. Feng, L. Ke, and K. Terada, "A GA-ACO hybrid algorithm for the multi-UAV mission planning problem," in *2014 14th International Symposium on Communications and Information Technologies (ISCIT)*, Sep. 2014, pp. 243–248.
- [34] Y. Chen, D. Yang, and J. Yu, "Multi-UAV Task Assignment With Parameter and Time-Sensitive Uncertainties Using Modified Two-Part Wolf Pack Search Algorithm," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 54, no. 6, pp. 2853–2872, Dec. 2018.
- [35] S. Gupta, G. Singal, and D. Garg, "Deep Reinforcement Learning Techniques in Diversified Domains: A Survey," *Archives of Computational Methods in Engineering*, vol. 28, no. 7, pp. 4715–4754, Dec. 2021.
- [36] X. Zhao, Z. Wang, and G. Zheng, "Two-Phase Neural Combinatorial Optimization with Reinforcement Learning for Agile Satellite Scheduling," *Journal of Aerospace Information Systems*, vol. 17, no. 7, pp. 346–357, 2020.
- [37] Y. Wang, S. Sun, and W. Li, "Hierarchical Reinforcement Learning for Vehicle Routing Problems with Time Windows," *Proceedings of the Canadian Conference on Artificial Intelligence*, Jun. 2021.
- [38] I. Bello, H. Pham, Q. Le, M. Norouzi, and S. Bengio, "Neural Combinatorial Optimization with Reinforcement Learning," Feb. 2017.
- [39] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer Networks," in *Advances in Neural Information Processing Systems*, vol. 28. Curran Associates, Inc., 2015.
- [40] Q. Ma, S. Ge, D. He, D. Thaker, and I. Drori, "Combinatorial Optimization by Graph Pointer Networks and Hierarchical Reinforcement Learning," *arXiv:1911.04936 [cs, stat]*, Nov. 2019.
- [41] Y. Xu, M. Fang, L. Chen, G. Xu, Y. Du, and C. Zhang, "Reinforcement Learning With Multiple Relational Attention for Solving Vehicle Routing Problems," *IEEE Transactions on Cybernetics*, pp. 1–14, 2021.
- [42] Y. Hu, Y. Yao, and W. S. Lee, "A reinforcement learning approach for optimizing multiple traveling salesman problems over graphs," *Knowledge-Based Systems*, vol. 204, p. 106244, Sep. 2020.
- [43] W. Liu, T. Zhang, R. Wang, K. Li, W. Li, and K. Yang, "Deep Reinforcement Learning for Orienteering Problems Based on Decomposition," Apr. 2022.
- [44] X. Xu, H. Yuan, M. Liptrott, and M. Trovati, "Two phase heuristic algorithm for the multiple-travelling salesman problem," *Soft Computing*, vol. 22, no. 19, pp. 6567–6581, Oct. 2018.
- [45] I. Dolinskaya, Z. E. Shi, and K. Smilowitz, "Adaptive orienteering problem with stochastic travel times," *Transportation Research Part E: Logistics and Transportation Review*, vol. 109, pp. 1–19, Jan. 2018.
- [46] I. V. Tetko, P. Karpov, R. Van Deursen, and G. Godin, "State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis," *Nature Communications*, vol. 11, no. 1, p. 5575, Nov. 2020.
- [47] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in *International Conference on Learning Representations*, Sep. 2020.
- [48] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, and P. Jiang, "BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. New York, NY, USA: Association for Computing Machinery, Nov. 2019, pp. 1441–1450.
- [49] J. Li, Y. Ma, R. Gao, Z. Cao, A. Lim, W. Song, and J. Zhang, "Deep Reinforcement Learning for Solving the Heterogeneous Capacitated Vehicle Routing Problem," *IEEE Transactions on Cybernetics*, pp. 1–14, 2021.
- [50] M. Nazari, A. Oroojlooy, L. Snyder, and M. Takac, "Reinforcement Learning for Solving the Vehicle Routing Problem," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.