

# A Two-stage Framework and Reinforcement Learning-based Optimization Algorithms for Complex Scheduling Problems

Yongming He, Guohua Wu, Yingwu Chen, and Witold Pedrycz, *Fellow, IEEE*

**Abstract**—There hardly exists a general solver that is efficient for scheduling problems due to their diversity and complexity. In this study, we develop a two-stage framework, in which reinforcement learning (RL) and traditional operations research (OR) algorithms are combined together to efficiently deal with complex scheduling problems. The scheduling problem is solved in two stages, including a finite Markov decision process (MDP) and a mixed-integer programming process, respectively. This offers a novel and general paradigm that combines RL with OR approaches to solving scheduling problems, which leverages the respective strengths of RL and OR: The MDP narrows down the search space of the original problem through an RL method, while the mixed-integer programming process is settled by an OR algorithm. These two stages are performed iteratively and interactively until the termination criterion has been met. Under this idea, two implementation versions of the combination methods of RL and OR are put forward. The agile Earth observation satellite scheduling problem is selected as an example to demonstrate the effectiveness of the proposed scheduling framework and methods. The convergence and generalization capability of the methods are verified by the performance of training scenarios, while the efficiency and accuracy are tested in 50 untrained scenarios. The results show that the proposed algorithms could stably and efficiently obtain satisfactory scheduling schemes for agile Earth observation satellite scheduling problems. In addition, it can be found that RL-based optimization algorithms have stronger scalability than non-learning algorithms. This work reveals the advantage of combining reinforcement learning methods with heuristic methods or mathematical programming methods for solving complex combinatorial optimization problems.

**Index Terms**—Scheduling, operations research, reinforcement learning, two-stage optimization.

## I. INTRODUCTION

**S**CHEDULING problems are widespread in production, manufacturing, transportation, and logistics, which are typically formulated as combinatorial optimization models. This type of model generally contains an objective function and a set of constraints, such as job-shop scheduling

problem (JSP)[1], traveling salesman problem with time window (TSPTW)[2], and Earth observation satellite scheduling problems[3]. Mathematical optimization models such as mixed-integer linear programming and nonlinear integer programming are popular for depicting concrete scheduling problems[4], but it is hard to find a relatively general and efficient algorithm to solve all of them, due to the significant difference in optimization objectives and the constraints of miscellaneous scheduling problems.

Tremendous number of scheduling problems in applications are proved to be NP-hard. In this study, a type of complex and representative problem is particularly concerned, which contains two levels of decision variables, i.e., 1) Which resource these tasks should be executed? 2) When should each candidate task execute?

For a long time, mathematical programming, heuristic and metaheuristic algorithms were the main tools to tackle the above scheduling problems[5]. However, the bottlenecks of these algorithms become prominent with the scale and the complexity of the problem increase: either the computing overhead is unacceptable, or the solution is unsatisfactory or unstable[3], [6]. Furthermore, the construction of rules or parameters in these algorithms is closely related to the characteristics of the problem, which leads to the fact that these algorithms are not versatile. Thus, a “general-purpose” and “far-sighted” algorithm with short decision time is needed.

The scheduling problem considered in this study actually can be decomposed into following three sub-problems[7]: 1) Assignment problem, for determining the possible tasks for each resource, i.e. assign tasks to different resources; 2) Sequencing problem, for determining the order of these tasks at each resource; 3) Timing problem, for determining the execution start and end time of scheduled tasks at each resource. These three sub-problems are mutually correlated, and the original scheduling problem could be solved by addressing them reasonably. It is a promising way to divide the original scheduling problems into two or multiple stages and then solve them separately, which could reduce the overall complexity and solve the scheduling problem more efficiently[8]. In particular, a two-stage framework is built for dealing with the considered scheduling problems.

The proposed two-stage framework consists of a prior stage and a rear stage: the prior stage for assignment problem, while the rear stage for sequencing problem and timing problem. It is desirable and natural to model the prior stage as a finite Markov decision process (MDP) and then solving it by

Yongming He is with the College of Systems Engineering, National University of Defense Technology, Changsha 410073, P. R. China. E-mail: heyongming10@hotmail.com

Guohua Wu (*the corresponding author*) is with the School of Traffic and Transportation Engineering, Central South University, Changsha 410075, China. E-mail: guohuawu@csu.edu.cn

Yingwu Chen is with the College of Systems Engineering, National University of Defense Technology, Changsha 410073, P. R. China. E-mail: ywchen@nudt.edu.cn

Witold Pedrycz is with the Department of Electrical and Computer Engineering, University of Alberta, Edmonton, AB T6G 2V4, Canada, and also with the Systems Research Institute, Polish Academy of Sciences, Warsaw 01447, Poland. E-mail: wpedrycz@ualberta.ca

reinforcement learning (RL): 1) The model can be well trained in advance by RL and achieve fast and accurate decisions in unknown scenarios; 2) It learns scheduling policy by trading off between immediate rewards and delayed value, and therefore RL makes decisions by weighing the estimated returns of each part of the system; 3) It is not required to manually design the decision-making rules for each instance. By contrast, the sub-problem in the rear stage can be formulated as a mixed-integer programming model, and then be solved by traditional operations research (OR) techniques.

For examining the effectiveness of our approach, deep Q-learning and an OR algorithm (i.e., a constructive heuristic or dynamic programming) are integrated the two-stage framework to solve the considered complex scheduling problem. The agile satellite scheduling problem is taken as an example to demonstrate the potential in applying the two-stage framework and the RL-based optimization algorithms. Experiments show that the two-stage scheduling framework is efficient in dealing with complex scheduling problems in the real-world.

The major contributions of this paper are:

- A novel two-stage scheduling framework is proposed for decomposing a type of complex scheduling problem into two parts, which are formulated by an MDP model and a mixed-integer programming model, respectively.
- An RL-based optimization mechanism is proposed under the two-stage framework. Specifically, deep Q-learning and OR algorithms (i.e., a constructive heuristic algorithm or a dynamic programming (DP) algorithm) are employed to solve the MDP model and mixed-integer programming model, respectively.
- The effectiveness of the proposed two-stage framework and RL-based optimization algorithms are verified by applying them to agile Earth observation satellite scheduling problem. In addition, some general conclusions about the model, mechanism, and algorithms are summarized.

The rest of the paper is structured as follows: Section II reviews existing studies on two-stage optimization models, scheduling algorithms, as well as reinforcement learning applications in combinatorial optimization problems. A general two-stage framework for scheduling problems is proposed in Section III. RL-based optimization methods are detailed in Section IV. Section V discusses the convergence and the generalization capabilities of proposed methods and then compares the proposed methods with other algorithms.

## II. LITERATURE REVIEW

As an important branch in the field of operations research, scheduling problems have attracted much attention from the scholars for a long time. After decades of research, scholars have been working to find more accurate and faster approaches to solve them, and some significant scheduling problems have been formed: resource-constrained project scheduling problem[9], vehicle routing problem with time windows (VRPTW)[10], job-shop scheduling problem (JSP)[11], and so on. Several ways for formulating these problems were proposed, such as the mathematical programming models[12], [13] and the constraint satisfaction model[14].

With the increasing real-life considerations in scheduling problems, the above models usually are hard to describe and simplify as constraints vary over different situations[15], [16]. Philippe Baptiste et al.[17] tried to summarize the constraints of scheduling problems into several categories, but many constraints in practical scheduling processes are still difficult to handle.

Because of the difficulty of many real-world scheduling problems, two-stage scheduling methods came into being. Two-stage decision-making architecture is common-used in practical scheduling projects, e.g., staff scheduling[18], [19], satellite task scheduling[20], issues in the transportation field[21], and extensive optimization problems in industry[22], [23]. However, solutions to these works are usually application-oriented. Hence, designing a novel and generic two-stage framework is meaningful to study scheduling problems in a better way.

It is of great significance to make satisfactory decisions with short computing time, especially when responding to scenario changes. Traditional algorithms for scheduling problems can be summarised into three categories: mathematical programming, heuristic algorithms, and metaheuristic algorithms. Mathematical programming such as branch-and-bound[24], branch-and-price[25], and dynamic programming[26] guarantee the optimal solution under certain assumptions, but the time complexity of these algorithms is usually exponential. Therefore, they are usually time-consuming and space-consuming in dealing with large scale problems. Heuristic algorithms find a solution by constructive rules for decision making. These rules are usually needed to be designed by experts based on experience. Heuristic algorithms can obtain a feasible solution in a short time. However, they cannot guarantee optimality, and may perform unsatisfactory in some situations. Furthermore, the search strategies of both exact algorithms and heuristic algorithms are highly coupled to the features and conditions of the concrete model considered. Metaheuristics such as tabu search[27], genetic algorithm[18], adaptive large neighborhood search (ALNS)[28] and their variants spring up due to their competitive performance in tackling complex scheduling problems. However, they are difficult to improve time efficiency and solution accuracy simultaneously, thereby the solutions may be unsatisfactory and unstable in affordable computation time.

It is no longer new about solving scheduling problems based on RL: Luca et.al.[1] proposed ant-Q algorithm for solving traveler salesman problem (TSP); Wei and Zhao[29] used Q-learning to select composite rule (a machine-job pair) in job-shop problem (JSP); Khalil et al.[30] suggested deep reinforcement learning (GNN-based deep Q-learning) to train the value function of each step of decision in the TSP problem; Nazari et.al.[31] solved vehicle routing problem (VRP) by pointer networks; Li et.al.[32] applied an improved actor-critic algorithm to train the model for multi-objective travelling salesman problem. These works show that RL has great potential in solving scheduling problems. However, the required number of iterations to reach satisfactory results rises sharply when RL is applied to intricate problems[6], which may be unacceptable in many real-world applications. Decomposing

the original problem and solving some sub-problems by OR algorithms can reduce the solution search space and improve learning speed of RL.

### III. THE TWO-STAGE SCHEDULING FRAMEWORK

For the sake of uniform and clear description of the type of scheduling problems considered in this study, we first introduce and discuss its characteristics and then build a two-stage scheduling framework for solving the problem efficiently.

#### A. Notation

Table I lists all necessary notation required to describe the scheduling problem, while algorithm-related notation will be illustrated where they occur.

TABLE I  
NOTATION

| Notation           | Interpretation  |
|--------------------|---|
| $i$                | Task index  |
| $j$                | Resource index  |
| $n$                | The number of tasks   |
| $m$                | The number of resources                                     |
| $TS$               | The set of tasks  |
| $RS$               | The set of resources  |
| $C$                | The capacity of resources                                   |
| $R_i$              | Reward at time step $i$                                     |
| $[ws_i^j, we_i^j]$ | Executable time window of task $i$ on resource $j$          |
| $d_i^j$            | Execution duration of task $i$ on resource $j$              |
| $p_i^j$            | Profit of task $i$ on resource $j$                          |
| $\Omega$           | Solution space of the original problem                      |
| $\Omega_1$         | Solution space of the problem in the prior stage            |
| $\Omega_2$         | Solution space of the problem in the rear stage             |
| $t$                | Time step in decision process                               |
| $S_t$              | State at time step $t$ in the MDP                           |
| $s_i^t$            | State of task $i$ at time step $t$ in the MDP               |
| $R_t$              | Reward at time step $t$ in the MDP                          |
| $\mathbf{A}$       | Action space of the MDP                                     |
| $\mathbf{A}(S_t)$  | Available action set under state $S_t$                      |
| $w_i^t$            | Remain executable time windows of task $i$ at time step $t$ |

#### B. Problem description

The concerned scheduling problems can be defined uniformly as follows[5]: There are a certain number of resources, each with certain capacities and restrictions in their usages. There are a certain number of tasks, each of which could associated with several records of executable time windows (ETWs), execution durations on each resource, and the profit of being completed. The problem is to assign the resource for scheduled tasks and then determine their execution start and end time, to optimize the objective function and meet all constraints. A schematic diagram of the problem is shown in Fig 1.

The input of a scheduling problem usually consists of two parts: resource information and task information.

Resource information describes the set of resources. Each resource may have capacity restrictions in multiple dimensions, e.g., The total amount of cargo that can be loaded on each vehicle in the VRP problem is limited by both volume

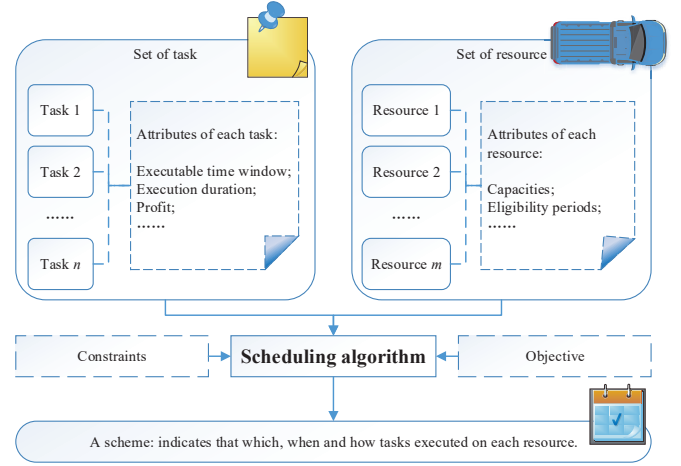


Fig. 1. A brief description of scheduling problems.

and weight. The set of resources information is formalized as follow:

$$RS = \{RS_1, RS_2, \dots, RS_m\} \quad (1)$$

$$RS_j = (C_1, C_2, \dots) \quad (2)$$

A task  $TS_i$  can be described by a set of triples: the ETW  $[ws_i^j, we_i^j]$ , required execution duration  $d_i^j$ , and the profit  $p_i^j$ . These attributes are related to the task  $i$  and resource  $j$ . The set of tasks is formalized as follow:

$$TS_i^j = \begin{cases} ([ws_i^j, we_i^j], d_i^j, p_i^j) & , i \text{ is available in } j \\ \emptyset & , \text{otherwise} \end{cases} \quad (3)$$

$$TS_i = \cup_{j=1,2,\dots,m} TS_i^j \quad (4)$$

$$TS = \cup_{i=1,2,\dots,n} TS_i \quad (5)$$

In general, the execution start time and the service resource are the decision variables for each scheduled task in the considered problem. They are represented by  $es_i$  and  $r_i$ , respectively.

$r_i$  is a variable indicating the resource on which task  $i$  is scheduled.

$$r_i = \begin{cases} -1, & \text{task } i \text{ is unscheduled} \\ j, & \text{task } i \text{ is scheduled on resource } j \end{cases} \quad (6)$$

$es_i$  and  $ee_i$  indicate the execution start and end time of task  $i$ , and the values of these two variables are *null* if and only if  $r_i = -1$ . Otherwise, we have

$$ee_i = es_i + d_i^{r_i}, r_i \neq -1 \quad (7)$$

If  $es_i \neq null$  and  $ee_i \neq null$ , the relationship between  $ws_i^j$ ,  $we_i^j$ ,  $es_i$  and  $ee_i$  is:

$$ws_i^j \leq es_i < ee_i \leq we_i^j, \exists j \in RS \quad (8)$$

### C. Two-stage representation for scheduling problems

Let us refer to the definition of two-stage optimization problem[33], [34], [35]. The two-stage scheduling model can be described as (9):

$$\begin{aligned} \min_{\mathbf{r}, \mathbf{es}} \quad & F(\mathbf{r}, \mathbf{es}) \\ \text{s.t.} \quad & G_1(\mathbf{r}) \leq \mathbf{0} \\ & G_2(\mathbf{r}, \mathbf{es}) \leq \mathbf{0} \\ & \mathbf{r} = (r_1, r_2, \dots, r_n) \in \Omega_1 \\ & \mathbf{es} = (es_1, es_2, \dots, es_n) \in \Omega_2 \end{aligned} \quad (9)$$

where,  $F(\mathbf{r}, \mathbf{es})$  is the objective function of the model;  $G_1(\mathbf{r})$  and  $G_2(\mathbf{r}, \mathbf{es})$  refers the constraints of the prior stage and the rear stage;  $\mathbf{r}$  and  $\mathbf{es}$  are decision variables;  $\Omega_1$  and  $\Omega_2$  denote the research spaces of the prior stage and the rear stage, respectively. The solution space of the original problem  $\Omega$  is calculated by the vector product of the solution space of two stages:  $\Omega = \Omega_1 \times \Omega_2$ . The feasible set of the rear stage could be reduced by fixing  $\mathbf{r}$ .

$$\Omega_2(\mathbf{r}') = \{\mathbf{es} \mid G_2(\mathbf{r}', \mathbf{es}) \leq \mathbf{0}\} \quad (10)$$

where,  $\mathbf{r}'$  is a fixed value of  $\mathbf{r}$ .

Generally, the rear stage can be formed as an optimization problem individually and solved by some optimization methods, as the problem of the rear stage is described as follows.

$$\begin{aligned} \min_{\mathbf{es} \in \Omega_2(\mathbf{r}')} \quad & F(\mathbf{r}', \mathbf{es}) \\ \text{s.t.} \quad & G_2(\mathbf{r}', \mathbf{es}) \leq \mathbf{0} \end{aligned} \quad (11)$$

However, the prior stage is hard to deal with in this way. The objective of the prior stage is the same as the objective of the original problem, and the objective value is strongly coupling with the results on the rear stage. In this study, we model the problem in the prior stage as an MDP model.

1) *Action*: A solution of the prior stage  $\mathbf{r}$  can be transformed to  $\mathcal{A}_t$ , which is a set of actions before any time step  $t$ , i.e.

$$\mathcal{A}_t = \{a_0, a_1, \dots, a_t\}, a_i \in \mathcal{A}(S_i), i = 0, 1, \dots, t \quad (12)$$

Actions in this model contain ‘‘Selecting a task for a resource’’ and ‘‘Termination’’. ‘‘Selecting a task for a resource’’ is the basic action in this model. Before taking an action at a time step, constraints which are unrelated to decision variables are checked before the task is selected for reducing unnecessary calculations. These constraints may include but are not limited to:

- Executable time window restriction;
- Consumption of each task and the capacity of the resource;

‘‘Termination’’ is an action that means the task selection within a certain planning cycle of this resource is completed. In the case of being without any available task at time step  $t$ , ‘‘Termination’’ is the only action.

2) *State*: A state depicts the attributes of each task on a determined time point of decision-making, it forms as Eq. (13):

$$S_t = \{s_i^t = (w_i^t, d_i^t, p_i^t) \mid t = 0, 1, \dots, n\} \quad (13)$$

where  $w_i^t$  is the set of remain ETWs of task  $i$  at time step  $t$ ;  $d_i^t$  represents the duration of task  $i$  at time step  $t$ ;  $p_i^t$  stands for the profit of completing task  $i$  at time step  $t$ .

The terminal states of scheduling problems are commonly expressed as Eq. (14), that is, there are no ETW remaining for all tasks.

$$w_i^t = \emptyset, \forall i \in \{1, 2, \dots, n\} \quad (14)$$

3) *Reward*: Reward indicates the difference in the value of the objective function between a step and its previous step. It can be counted by Eq. (15)

$$R_t = F(\mathcal{A}_t, \mathbf{es}) - F(\mathcal{A}_{t-1}, \mathbf{es}) \quad (15)$$

where  $F(\mathcal{A}_t, \mathbf{es})$  and  $F(\mathcal{A}_{t-1}, \mathbf{es})$  are the values of objective function at time step  $t$  and  $t - 1$ , respectively.

4) *Value function*: Value function refers to the expected long term profit under the condition of determined state and action, which is the criterion of taking actions. Due to the huge size of the state, tabular solution methods are hard to reflect the nature of the value function and support decisions efficiently. Artificial neural networks (ANNs) are popular and efficient to represent the value function owing to their strong prediction and classification capabilities[36], [37], so we deem that ANN could be a promising tool for estimating the value function.

The input of the network is the state of the model, while the output is the expected long term profit of each action. The value function is continuously updated according to a large number of records about actions, states, and rewards.

## IV. REINFORCEMENT LEARNING-BASED OPTIMIZATION ALGORITHMS

We decompose the scheduling problem into two stages, i.e., a finite Markov decision process for assignment problem while a mixed-integer programming process for sequencing problem and timing problem. The interactions between these two processes are that mixed-integer programming process obtains a solution according to the results of finite MDP; while the finite MDP makes decisions considering the results of the mixed-integer programming. During training phase, the value function of the MDP is continually fitted based on records of decisions and results, for a better decision criterion. The relationship between the two stages is shown in Fig. 2.

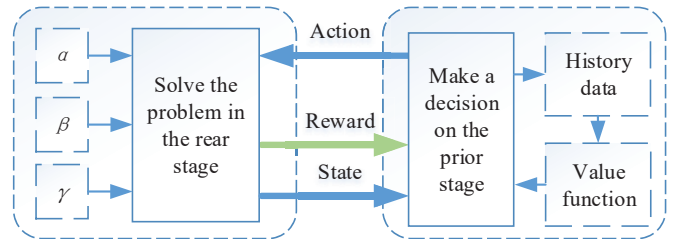


Fig. 2. The interactions between two stages of the process.  $\alpha, \beta$ , and  $\gamma$  draw the environment, objective, and constraints of the mixed-integer programming problem, respectively; Action, Reward, and State are defined in Section III.

We design the deep Q-learning (DQN) for the assignment problem in the prior stage, and an OR algorithm (a constructive heuristic (HADRT) or dynamic programming(DP)) for the sequencing and timing problem in the rear stage. The approaches with HADRT and DP are named DQN\_CH and DQN\_DP, respectively.

### A. DQN in the prior stage

There is a wide range of state of scheduling problems. In the other hand, the initial states of a scheduling problem vary because of the changable attributes of tasks and resources, so it is necessary to get adequate knowledge in limited iterations and to make decisions under unknown states.

The training process is usually organized as the following steps:

**Step 1:** Choose an action by “Exploitation” or “Exploration”;

**Step 2:** If the action is not “Termination”, get the reward and state by interactions between two stages;

**Step 3:** Update the value function based on actions, rewards, and states;

**Step 4:** Repeat **Step 1** to **Step 3** until the termination state is reached.

The policy for choosing actions includes two strategies, i.e., “Exploitation” and “Exploration”. “Exploitation” signifies choosing the best action according to the predicted value function, while “Exploration” designates selecting other actions randomly.

The value function is updated as Eq. (16)

$$\hat{q}(S_t, A_t) = R_{t+1} + \gamma \max_A q(S_{t+1}, A_{t+1}) \quad (16)$$

where  $\hat{q}$  is the target Q function;  $q$  is the estimated Q function;  $R_{t+1}$  is the reward of time step  $t+1$ ;  $\gamma$  is a constant to indicate the discount rate.

No prior knowledge about the value function required at this process, instead, a more accurate value function can be fitted through experiences. The pseudo-code of choosing an action is as follows:

---

#### Algorithm 1: Choosing a task by “Exploitation” or “Exploration”

---

**Input:** State at time step  $t$   $S_t$ , available tasks at time step  $t$   $A(S_t)$ , and the Q-network  $Q$

**Output:** A selected task from all available tasks  $a_t$

```

1 Threshold  $\varepsilon$  initialization
2  $idx \leftarrow \text{random}()$ 
3 if  $idx \leq \varepsilon$  then
4    $a_t \leftarrow$  A task which is randomly chosen from  $AL_t$ 
5 else
6   Calculate  $Q(S_t)$ 
7    $Q\_in\_list \leftarrow \emptyset$ 
8   for  $al \in A(S_t)$  do
9      $Q\_in\_list \leftarrow Q\_in\_list.append(Q(s, AL[i]))$ 
10   $a_t \leftarrow \text{argmax}(Q\_in\_list)$ 

```

---

In the above algorithm, threshold value  $\varepsilon$  is the probability of choosing an action randomly (line 3 to 4). Otherwise, the action is taken according to  $Q$ . The Q values of all actions at state  $s_t$  are calculated by the Q-network, which is trained by Eq. (16). Only the actions in  $A(S_t)$  are considered at time step  $t$ , so we record the Q value of all available actions  $Q\_in\_list$ , and then get the action with the highest Q value. see line 7 to 10.

After training, the value function (Q-network) has been finalized. In the testing process, the algorithm always select the action with the highest Q value for seeking the satisfactory solution.

### B. OR algorithms in the rear stage

Once the solution of the prior stage has been determined, there are several types of algorithms for dealing with the problem in the rear stage. Dynamic programming and a constructive heuristic algorithm are applied in this study.

1) *Constructive heuristic algorithm:* A constructive heuristic algorithm based on the density of residual tasks, named HADRT[38], is introduced as one of the algorithms for sequencing and timing problem. The core of HADRT are the expressions below:

$$f(i) = g(i) + h(i) = \sum_{i|r_i>0} 1 + \frac{ET - we_i}{\bar{d}_i} \quad (17)$$

$$k = \arg \max f(i) \quad (18)$$

where,  $g(i)$  is the number of scheduled tasks after adding the  $i$ -th task;  $h(i)$  equals the ratio between the remaining time length and the average duration of tasks.

Under the premise of satisfying all constraints, execution start time for each scheduled task is determined as follows:

$$es_k = \max(wb_k, ee_{k-1} + ct_{k,k-1}) \quad (19)$$

where  $ct_{k,k-1}$  indicates the preparation time required from the end of task  $k-1$  to the start of task  $k$ .

The expectation of total profit in the scheme obtained by HADRT is maximized, under the following conditions:

- Values of preparation time and durations of tasks are fixed or follow a determined distribution;
- The profit of each task is a constant;
- The ETW of any task is not completely covered by the ETW of another task.

The above statement has been proven by mathematical induction and reduction. The detailed process of the proof is derived in [38], which is not discussed here. The algorithm is organized as the sequence of the following steps:

**Step 1:** Initialization: the latest end time of the scenario, the set of candidate tasks;

**Step 2:** Calculate  $f(i)$  for each available task;

**Step 3:** Find the task  $k$  by Eq. (17) and (18);

**Step 4:** Set the execution start time of task  $k$  by Eq. (19);

**Step 5:** Update the latest end time of the scenario if the task  $k$  is scheduled;

**Step 6:** Repeat **Step 2** to **Step 5** until all of the tasks are considered.

2) *Dynamic programming:* Dynamic programming (DP) is one of the commonly used methods for solving combinatorial optimization problems. DP guarantees the optimal solution for the timing problem in the rear stage with time complexity  $O(n^2)$ , under the assumptions of: 1) The sequence of tasks is fixed. 2) Time could be discretized. The above statement can be proved because the problem meets the optimality principle of dynamic programming[39].

The sequence for timing problem can be obtained by a greedy algorithm, that is, the tasks are sorted in ascending order of the end time of ETWs.

The schematic way of applying DP to this problem is displayed in Fig. 3.

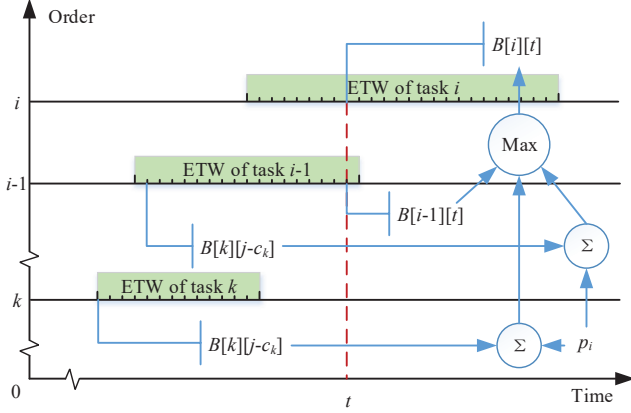


Fig. 3. The schematic diagram of dynamic programming in the problem.

It is worth noting that this process meets the principle of optimality and without aftereffect[39]. As shown in Figure 3,  $B[i][t]$  is the maximum total profit that the top  $i$ -th tasks can get in the time  $t$ , which is calculated as

$$B[i][t] = \max_k \{B[i-1][t], B[i-k][t - d_i - ct_{i,k}] + p_i\} \quad (20)$$

Based on the above, the pseudo-code of solving the timing problem by DP is as follows:

---

**Algorithm 2: Dynamic programming**

---

**Input:** A sequence of tasks  $T$

**Output:** Objective function value of the optimal solution  $B[|T|][MaxTime]$

```

1 Parameters initialization
2  $MaxTime \leftarrow$  timespan
3 foreach  $i \in T$  do
4   for  $t = 0 : MaxTime$  do
5     for  $k = 0 : i$  do
6        $TempB \leftarrow B[i-k][t - d_i - ct_{i,k}] + p_i$ 
7       if satisfy all constraints and
8          $TempB > B[i][t]$  then
9            $B[i][t] \leftarrow TempB$ 
10        else
11           $B[i][t] \leftarrow B[i-1][t]$ 

```

---

3) *Summary on HADRT and DP:* HADRT and DP are representative heuristic algorithm and mathematical programming algorithm in traditional OR area, respectively. Reward of time step  $t$  is calculated by the result of HADRT or DP at time step  $t$ , therefore, the performance of HADRT and DP affects the rewards of RL. The range of rewards obtained by HADRT should be contained in the interval

$[-\max_{i,j} p_i^j, \max_{i,j} p_i^j]$ , in which  $p_i^j$  is the profit of task  $i$  on resource  $j$ . Correspondingly, rewards obtained by DP should be non-negative because DP guarantees the optimality of timing problem, and the range of rewards obtained by DP is reduced to  $[0, \max_{i,j} p_i^j]$ .

Both the time complexity of HADRT and DP are  $O(n^2)$ , in which,  $n$  represents the scale of the problem. Therefore, the time complexity of these two algorithms is in a low level.

We discuss the performance and consumption of these two approaches in the experiments.

## V. EXPERIMENTAL STUDIES

Experiments are implemented on a personal computer with Intel(R) Core(TM) i7-8750H CPU, 16.0GB RAM and NVIDIA GeForce GTX 1060. The agile Earth observation satellite (AEOS) scheduling problem is selected as a representative example to verify the model and algorithm proposed in this paper. This problem is recognized as one of the complicated practical problems in the scheduling field[40].

### A. Definition of AEOS scheduling problem

This section defines AEOS scheduling problem from inputs, objective function, and constraints.

1) *Design of resources:* Due to the periodicity of the satellite orbit, the scheduling process of the satellite can be naturally divided into orbital periods. Thereby, available orbital periods of every satellite are viewed as resources in AEOS scheduling problem.

Refer to Eq. (2), an orbital resource in AEOS scheduling problem could be described as the tuple as Eq. (21):

$$RS_j = (Energy_j, Storage_j) \quad (21)$$

A constant  $Energy_j$  shows the available energy of resource  $j$  and a constant  $Storage_j$  indicates the available storage. In our simulation, the values of  $Energy_j$  and  $Storage_j$  are 150 Ah and 2000 GB for all resources, which are consistent with the standards of the satellite industry.

2) *Design of tasks:* A task in AEOS scheduling problem could be described by referring to Eq. (3) and Eq. (4). Under our settings,  $p_i^j$  and  $d_i^j$  are predetermined;  $ws_i^j$  and  $we_i^j$  involved in the time windows of tasks could be calculated prior to the scheduling process by considering the orbit parameters and task locations  $(lat_i, lon_i)$ . The discretization interval of  $ws_i^j$  and  $we_i^j$  are 1 second.

The diversity of AEOS scheduling problem leads to the difficulty of finding a generally recognized benchmark that is close to actual applications. To facilitate comparisons and analyses, we create several test sets following certain distributions. Rules for generating simulated tasks are listed in Table II.

Two types of scenarios were designed to test the effectiveness of the proposed method in different geographical distributions of tasks. The ability to handle routine tasks could be tested in large-area scenes, which covers the whole are of China; the responsiveness of AEOS in emergencies like floods and earthquakes could be tested in small-area scenes, which covers Hunan province in China. Accurately, two small-area

TABLE II  
RULES FOR GENERATING DETAILS OF TASKS

| Parameters                   | Distribution | Type    | Lower bound | Upper bound |
|------------------------------|--------------|---------|-------------|-------------|
| $lat_i$ in large-area scenes | Evenly       | Float   | 3           | 53          |
| $lon_i$ in large-area scenes | Evenly       | Float   | 73          | 133         |
| $lat_i$ in small-area scenes | Evenly       | Float   | 20          | 30          |
| $lon_i$ in small-area scenes | Evenly       | Float   | 108         | 114         |
| $d_i^j$                      | Constant     | Integer | 5           | 5           |
| $p_i^j$                      | Evenly       | Integer | 1           | 10          |

scenes with task sizes of 20 and 50 as well as three large-area scenes with task sizes of 100, 200, and 400 were designated, which are marked as H\_20, H\_50, C\_100, C\_200, and C\_400, respectively.

3) *Objective and constraints*: The objective function and constraints of AEOS scheduling models are usually depended on the specific satellite platforms and application mode of the observation systems. By extensively investigating the descriptions of AEOS scheduling problems in literatures, we summarize the common objective function and some common constraints of the problem[20], [28], [41], [42].

$$\max \sum_{\{i|r_i>0\}} p_i^{r_i} \quad (22)$$

subject to

$$\text{card}(\{r_i | r_i > 0\}) \leq 1 \quad (23)$$

$$\sum_{\{i|r_i=j\}} \text{Storage}(i, j) \leq \text{Maxstorage}(j) \quad (24)$$

$$\sum_{\{i|r_i=j\}} \text{Energy}(i, j) \leq \text{MaxEnergy}(j) \quad (25)$$

$$ws_i^{r_i} - es_i \leq 0 \quad (26)$$

$$es_i + d_i^{r_i} - we_i^{r_i} \leq 0 \quad (27)$$

$$(es_{i_0} - es_{i_1})(es_{i_0} - ee_{i_1}) > 0 \quad \forall i_0 \neq i_1, r_{i_0} = r_{i_1} \quad (28)$$

$$(ee_{i_0} - es_{i_1})(ee_{i_0} - ee_{i_1}) > 0 \quad \forall i_0 \neq i_1, r_{i_0} = r_{i_1} \quad (29)$$

$$|es_{i_0} - ee_{i_1}| \geq MST(i_0, i_1) \quad \forall i_0 \neq i_1, r_{i_0} = r_{i_1} \quad (30)$$

$$i = 1, 2, \dots, n \quad (31)$$

$$j = 1, 2, \dots, m \quad (32)$$

In the above model, the objective function is to maximize the sum of the profit of all scheduled tasks, which as (22). Inequality (23) to (25) affiliate to constraints in the prior stage, while (26) to (30) belong to that of the rear stage.

Constraints (23) refers that a task is scheduled at most once; (24) and (25) shows the capacity constraints of the AEOSs with respect to energy and storage. (26) and (27) indicate that all tasks should be executed in one of their ETWs. (28) to (30) expresses the constraints between two tasks, i.e., no overlap in the execution time of any two scheduled tasks is allowed. Meanwhile, the interval between the two tasks must be greater than the corresponding slew time, to guarantee the shortest time required by the satellite to adjust the attitude from a task to the next one.

The configurations of slew time, storage consumption, and electricity consumption are derived from a practical example in China, which named AS-01. The necessary functions of AS-01 are detailed in [28], [43].

### B. Performance of applying Reinforcement Learning-based algorithms to AEOS scheduling

The feasibilities of the proposed algorithms are tested in this section. In addition to DQN\_CH and DQN\_DP that are introduced in Section IV, the experiment introduces two other approaches for comparison of results, that is:

**DQN\_CH/C**: DQN for assignment problem and sequencing problem, while HADRT for timing problem;

**DQN\_DP/C**: DQN for assignment problem and sequencing problem, while DP for timing problem;

1) *Convergence analysis*: 1000 iterations on the same instance are used to verify the convergence of the proposed methods, as shown in Fig. 4. Through dozens of iterations, the total profit of all algorithms is stable overall, with acceptable fluctuations. The fluctuation is mainly caused by the existence of ‘‘Exploration’’ in the training process of DQN.

As the number of tasks increases, the differences between these four algorithms are apparent: The average total profit of DQN\_CH and DQN\_DP in all scenarios is higher than DQN\_DP/C and DQN\_CH/C. Furthermore, DQN\_DP/C and DQN\_CH/C will converge to a low level on C\_100, C\_200, and C\_400. These two algorithms are easily trapped in a local optimum solution, and they did not show the trend to escape the local optimum in the first 1000 iterations. The results show the sequencing problem is a high-complexity problem for DQN. Thus, these two algorithms are hard to achieve satisfactory results through small-scale training. The total profit of DQN\_DP averages around 10% higher than that of DQN\_CH, but DQN\_DP requires more computing resources and time during the training process: it can only complete 949 iterations in C\_400 because of memory overflow.

2) *Generalization analysis*: 20 random instances are used to test the generalization of the proposed methods. Depending on the different scales, training processes of different instances take tens of minutes to several hours, including the time for preprocessing and status updates. Since the profit of each task in a task set is random, analyzing the trend of profit margin is more reasonable. The profit margin is the ratio of the total profit of scheduled tasks to that of all candidate tasks. As can be seen in Fig. 5, although only 20 iterations are performed on each data set, the polyline of each algorithm fluctuates across a certain level, indicating the value functions trained by all four algorithms are applicable in different scenarios, that is, these four algorithms are robust in generalization. However, the fluctuations in Fig. 5 are larger than the fluctuation range in Fig. 4. This is because: 1) the optimal solutions vary in each instance; 2) the value functions need to be continuously modified in different instances.

The stability of the results in distinct instances implies the generalization. Algorithms trained from 20 instances are applied to unknown scenes, and results are shown in Fig. 6. In 50 instances with different features, the test results of the four

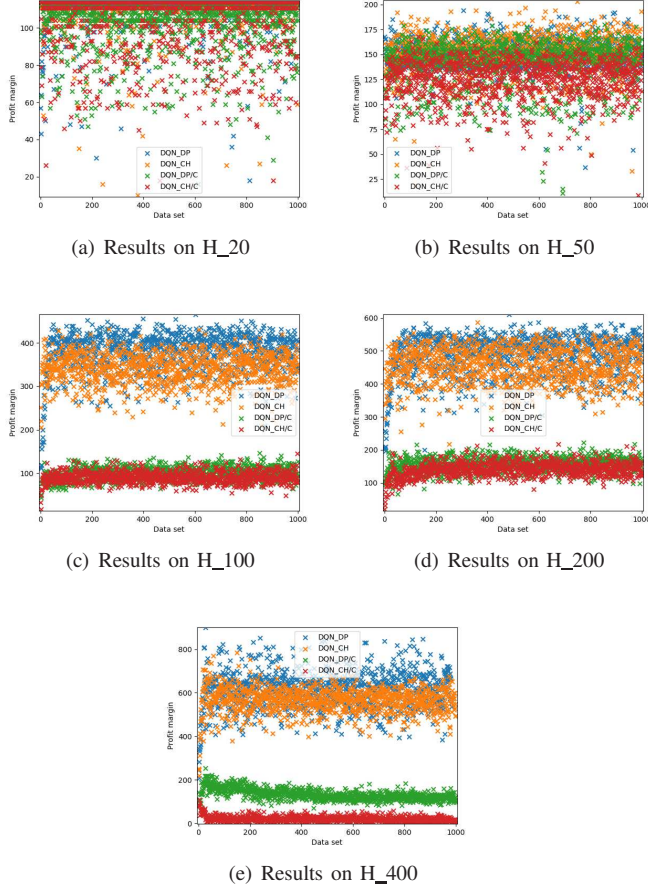


Fig. 4. Total profit of proposed algorithms under 1000 iterations in single instance.

algorithms are consistent with the rules during training. The results of DQN\_CH and DQN\_DP are much better than those of the other two. In all 50 test sets, DQN\_CH are superior to DQN\_DP on just two instances; except for the instances in H<sub>20</sub>, the result of DQN\_CH and DQN\_DP equally on only one instance.

### C. Comparison with non-learning scheduling algorithms

We compared DQN\_CH and DQN\_DP to several advanced algorithms that have published in recent papers. These algorithms include a heuristic represented by a heuristic algorithm based on the density of residual tasks[38], a meta-heuristic represented by the adaptive large neighborhood search algorithm[28], and a mathematical programming method represented by the branch and bound algorithm[24], which are denoted as HADRT, ALNS, and B&B, respectively.

1) *Total profit*: In this problem, total profit is the value of the objective function, which shows the optimization capability of an algorithm. Thus, it is a vital indicator for evaluating algorithms. We ran concerned algorithms in all test sets and recorded results whose running time less than 3600 seconds. All results are summarized in Fig. 7. The box plot shows: 1) B&B could obtain the optimal solution in small scale problems (in test set H<sub>20</sub>), but cannot obtain the result

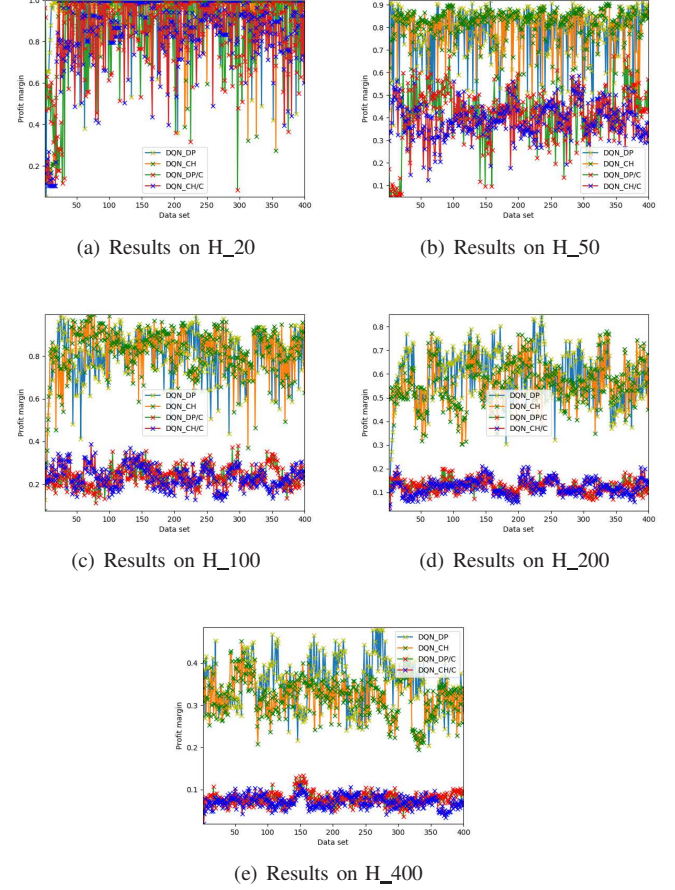


Fig. 5. Total profit of proposed algorithms under 20 iterations in multiple instances.

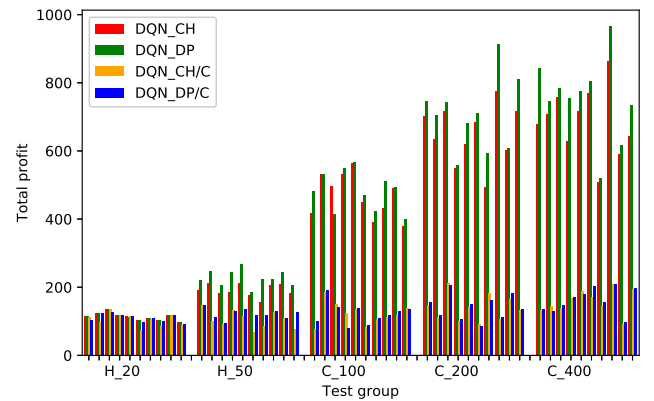


Fig. 6. Total profit of proposed algorithms on test data.

within 3600 seconds in other test sets; 2) HADRT could reach good results for most instances, but in some cases, the results are not satisfactory; 3) ALNS, DQN\_CH, and DQN\_DP can reliably achieve high-quality results, especially for large test sets. In general, DQN\_DP performs the best among these five comparative algorithms.

2) *Running time*: For practical problems, the running time of the program is another important metric to evaluate the



TABLE III  
THE AVERAGE AND STANDARD DEVIATION OF RUNNING TIME ON DIFFERENT DATA SETS.

| Scenes | MEAN   |        |       |        |       | STDEV  |        |       |        |       |
|--------|--------|--------|-------|--------|-------|--------|--------|-------|--------|-------|
|        | DQN_CH | DQN_DP | HADRT | ALNS   | B&B   | DQN_CH | DQN_DP | HADRT | ALNS   | B&B   |
| H_20   | 0.079  | 1.824  | 0.038 | 7.777  | 1.553 | 0.004  | 0.038  | 0.013 | 1.146  | 1.234 |
| H_50   | 0.153  | 5.539  | 0.126 | 20.287 | >3600 | 0.028  | 0.38   | 0.011 | 11.653 | —     |
| C_100  | 0.399  | 9.367  | 0.152 | 38.691 | >3600 | 0.024  | 0.752  | 0.017 | 26.916 | —     |
| C_200  | 1.003  | 24.655 | 0.442 | 316.83 | >3600 | 0.12   | 1.44   | 0.034 | 75.318 | —     |
| C_400  | 1.809  | 24.426 | 1.246 | 2057.9 | >3600 | 0.234  | 2.031  | 0.098 | 222.41 | —     |

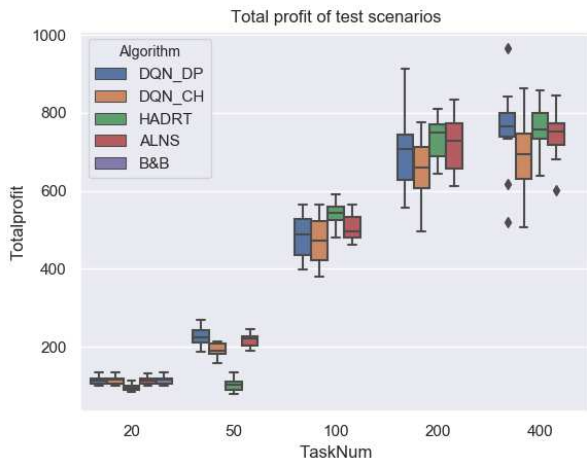


Fig. 7. Comparison of the RL-based optimization algorithms with advanced scheduling algorithms on the test data.

efficiency of the algorithm. According to the data in Table III, HADRT, DQN\_CH, and DQN\_DP could get results in a short time, while the running time of ALNS and B&B grow dramatically with the increase of the problem size. The computation time of DQN\_CH is approximately twice that of HADRT. Although the running time of DQN\_DP is much larger than that of DQN\_CH and HADRT, they increase linearly with the scale of the problem. By analyzing the standard deviation of running time, it can be found that under the same test group, the computation time of ALNS and B&B is greatly affected by the characteristics of inputs. By contrast, HADRT, DQN\_CH, and DQN\_DP run with a stable time.

Overall, B&B only works well on small-scale instances. HADRT could find a feasible solution in a short time, but the quality of the solution could not be guaranteed. For ALNS, the risk of falling into a local optimum can be reduced through the random search mechanism, but it usually takes too long to search a promising solution, which may be unacceptable in many practical problems. RL-based optimization algorithms could be a potential alternative to obtain satisfactory solutions with acceptable time.

## VI. CONCLUSIONS

We have summarized the essential characteristics of a type of complex scheduling problems, and then propose a two-stage scheduling framework to solve them. Based on the framework,

we have built two RL-based optimization algorithms. These algorithms are designed by combining DQN with different OR algorithms (i.e., HADRT or DP). Experimental results demonstrate the effectiveness and advantage of RL-based optimization algorithms.

This is a preliminary attempt about combining the machine learning and operations research to solve a practical scheduling problem. A two-stage decision model is devised to describe general scheduling problem. Decision variables are divided into two parts and their values are determined in two stages. The assignment problem can be seen as the prior stage and solved by RL, while the OR algorithms can be used to solve the sub-problems in the rear stage. It is more likely to fall into a local optimum when solving the sequencing problem in the prior stage, and it is difficult to jump out of the local optimum in small-scale training. The impact of different OR algorithms used in the rear stage is analyzed. The reward obtained by DP is closer to the optimal solution of the problem in the rear stage than that of HADRT, thereby DQN\_DP could perform better. However, DQN\_DP consumes relatively more computing resources both in training and testing.

Although deep Q-learning is one of the basic methods in RL, this work shows the great potential in solving complex combinatorial optimization problems by the RL-based optimization algorithms. To further enrich the conclusions, we will try more RL methods and structures of the value function, and delve into the inherent rules of the integrating learning algorithms with traditional optimization algorithms. Furthermore, we will continue to test the performance of RL-based optimization algorithms on other types of scheduling problems.

## VII. ACKNOWLEDGMENTS

This research is supported by the Natural Science Fund for Distinguished Young Scholars of Hunan Province under Grant 2019JJ20026 and the National Natural Science Foundation of China under Grant No.71701204.

## REFERENCES

- [1] L. M. Gambardella and M. Dorigo, "Ant-q: A reinforcement learning approach to the traveling salesman problem," in *Machine Learning Proceedings 1995*. Elsevier, 1995, pp. 252–260.
- [2] H. Kona, A. Burde, and D. Zanwar, "A review of traveling salesman problem with time window constraint," *IJIRST-International Journal for Innovative Research in Science & Technology*, vol. 2, 2015.
- [3] Y. Du, T. Wang, B. Xin, L. Wang, Y. Chen, and L. Xing, "A data-driven parallel scheduling approach for multiple agile earth observation satellites," *IEEE Transactions on Evolutionary Computation*, 2019.

- [4] B. Ji, X. Yuan, and Y. Yuan, "A hybrid intelligent approach for scheduling of cascaded locks with multiple chambers," *IEEE transactions on cybernetics*, vol. 49, no. 4, pp. 1236–1248, 2018.
- [5] L. P. Michael, *Scheduling: theory, algorithms, and systems*. Springer, 2018.
- [6] H. Lu, X. Zhang, and S. Yang, "A learning-based iterative method for solving vehicle routing problems," in *International Conference on Learning Representations*, 2019.
- [7] W. J. Wolfe and S. E. Sorensen, "Three scheduling algorithms applied to the earth observing systems domain," *Management Science*, vol. 46, no. 1, pp. 148–166, 2000.
- [8] Y. Wang, Z.-Y. Ru, K. Wang, and P.-Q. Huang, "Joint deployment and task scheduling optimization for large-scale mobile users in multi-uav-enabled mobile edge computing," *IEEE Transactions on Cybernetics*, 2019.
- [9] S. Hartmann and D. Briskorn, "A survey of variants and extensions of the resource-constrained project scheduling problem," *European Journal of operational research*, vol. 207, no. 1, pp. 1–14, 2010.
- [10] A. Dixit, A. Mishra, and A. Shukla, "Vehicle routing problem with time windows using meta-heuristic algorithms: a survey," in *Harmony Search and Nature Inspired Optimization Algorithms*. Springer, 2019, pp. 539–546.
- [11] A. Jones, L. C. Rabelo, and A. T. Sharawi, "Survey of job shop scheduling techniques," *Wiley encyclopedia of electrical and electronics engineering*, 2001.
- [12] J. Bracken and J. T. McGill, "Mathematical programs with optimization problems in the constraints," *Operations Research*, vol. 21, no. 1, pp. 37–44, 1973.
- [13] D. Bertsimas and J. N. Tsitsiklis, *Introduction to linear optimization*. Athena Scientific Belmont, MA, 1997, vol. 6.
- [14] G. Verfaillie, M. Lemaître, and T. Schiex, "Russian doll search for solving constraint optimization problems," in *AAAI/IAAI, Vol. 1*, 1996, pp. 181–187.
- [15] K. Lachhwani and A. Dwivedi, "Bi-level and multi-level programming problems: Taxonomy of literature review and research issues," *Archives of Computational Methods in Engineering*, vol. 25, no. 4, pp. 847–877, 2018.
- [16] J. Turner, Q. Meng, G. Schaefer, A. Whitbrook, and A. Soltoggio, "Distributed task rescheduling with time constraints for the optimization of total task allocations in a multirobot system," *IEEE transactions on cybernetics*, vol. 48, no. 9, pp. 2583–2597, 2017.
- [17] P. Baptiste, P. Laborie, C. Le Pape, and W. Nuijten, "Constraint-based scheduling and planning," in *Foundations of artificial intelligence*. Elsevier, 2006, vol. 2, pp. 761–799.
- [18] C.-C. Tsai and S. H. Li, "A two-stage modeling with genetic algorithms for the nurse scheduling problem," *Expert Systems with Applications*, vol. 36, no. 5, pp. 9506–9512, 2009.
- [19] A. Parisio and C. N. Jones, "A two-stage stochastic programming approach to employee scheduling in retail outlets with uncertain demand," *Omega*, vol. 53, pp. 97–103, 2015.
- [20] B. Deng, C. Jiang, L. Kuang, S. Guo, J. Lu, and S. Zhao, "Two-phase task scheduling in data relay satellite systems," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 2, pp. 1782–1793, 2017.
- [21] F. Wu and R. Sioshansi, "A two-stage stochastic optimization model for scheduling electric vehicle charging loads to relieve distribution-system constraints," *Transportation Research Part B: Methodological*, vol. 102, pp. 55–82, 2017.
- [22] H. Qiu, B. Zhao, W. Gu, and R. Bo, "Bi-level two-stage robust optimal scheduling for ac/dc hybrid multi-microgrids," *IEEE Transactions on Smart Grid*, vol. 9, no. 5, pp. 5455–5466, 2018.
- [23] H. Dasthi, A. J. Conejo, R. Jiang, and J. Wang, "Weekly two-stage robust generation scheduling for hydrothermal power systems," *IEEE Transactions on Power Systems*, vol. 31, no. 6, pp. 4554–4564, 2016.
- [24] X. Chu, Y. Chen, and L. Xing, "A branch and bound algorithm for agile earth observation satellite scheduling," *Discrete Dynamics in Nature and Society*, vol. 2017, 2017.
- [25] R. Václavík, A. Novák, P. Šoucha, and Z. Hanzálek, "Accelerating the branch-and-price algorithm using machine learning," *European Journal of Operational Research*, vol. 271, no. 3, pp. 1055–1069, 2018.
- [26] Y. Xiao and A. Konak, "A genetic algorithm with exact dynamic programming for the green vehicle routing & scheduling problem," *Journal of Cleaner Production*, vol. 167, pp. 1450–1463, 2017.
- [27] O. Bräysy and M. Gendreau, "Tabu search heuristics for the vehicle routing problem with time windows," *Top*, vol. 10, no. 2, pp. 211–237, 2002.
- [28] X. Liu, G. Laporte, Y. Chen, and R. He, "An adaptive large neighborhood search metaheuristic for agile satellite scheduling with time-dependent transition time," *Computers & Operations Research*, vol. 86, pp. 41–53, 2017.
- [29] Y. Wei and M. Zhao, "A reinforcement learning-based approach to dynamic job-shop scheduling," *Acta Automatica Sinica*, vol. 31, no. 5, p. 765, 2005.
- [30] E. Khalil, H. Dai, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," in *Advances in Neural Information Processing Systems*, 2017, pp. 6348–6358.
- [31] M. Nazari, A. Oroojlooy, L. Snyder, and M. Takác, "Reinforcement learning for solving the vehicle routing problem," in *Advances in Neural Information Processing Systems*, 2018, pp. 9839–9849.
- [32] K. Li, T. Zhang, and R. Wang, "Deep reinforcement learning for multiobjective optimization," *IEEE Transactions on Cybernetics*, 2020.
- [33] M. Goerigk, A. Kasperski, and P. Zieliński, "Two-stage combinatorial optimization problems under risk," *Theoretical Computer Science*, vol. 804, pp. 29–45, 2020.
- [34] L. N. Vicente and P. H. Calamai, "Bilevel and multilevel programming: A bibliography review," *Journal of Global optimization*, vol. 5, no. 3, pp. 291–306, 1994.
- [35] J. Lu, J. Han, Y. Hu, and G. Zhang, "Multilevel decision-making: A survey," *Information Sciences*, vol. 346, pp. 463–487, 2016.
- [36] Y. Bengio, A. Lodi, and A. Prouvost, "Machine learning for combinatorial optimization: a methodological tour d'horizon," *arXiv preprint arXiv:1811.06128*, 2018.
- [37] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multi-agent systems: A review of challenges," *Solutions and Applications*, 2018.
- [38] Y. He, Y. Chen, J. Lu, C. Chen, and G. Wu, "Scheduling multiple agile earth observation satellites with an edge computing framework and a constructive heuristic algorithm," *Journal of Systems Architecture*, vol. 95, pp. 55–66, 2019.
- [39] R. Bellman, "Dynamic programming," *Science*, vol. 153, no. 3731, pp. 34–37, 1966.
- [40] S. Mitrovic-Minic, D. Thomson, J. Berger, and J. Secker, "Collection planning and scheduling for multiple heterogeneous satellite missions: Survey, optimization problem, and mathematical programming formulation," in *Modeling and Optimization in Space Engineering*. Springer, 2019, pp. 271–305.
- [41] G. Wu, W. Pedrycz, H. Li, M. Ma, and J. Liu, "Coordinated planning of heterogeneous earth observation resources," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 46, no. 1, pp. 109–125, 2015.
- [42] K. Luo, H. Wang, Y. Li, and Q. Li, "High-performance technique for satellite range scheduling," *Computers & Operations Research*, vol. 85, pp. 12–21, 2017.
- [43] Y. He, "Research on agile satellite autonomous mission planning system and re-planning method," Master's thesis, National University of Defense Technology, 2016.



**Yongming He** received the B.S. degree in Logistics Engineering from Chang'an University, China, in 2014. He is now pursuing his Ph.D degree in Management Science and Engineering at the College of Systems Engineering, National University of Defense Technology, China. He is a visiting Ph.D. student at University of Alberta, Edmonton, AB, Canada, from Nov. 2018 to Nov. 2019. His research interests include operations research, artificial intelligence, intelligent decision, scheduling and planning.



**Guohua Wu** received the B.S. degree in Information Systems and Ph.D degree in Operations Research from National University of Defense Technology, China, in 2008 and 2014, respectively. During 2012 and 2014, he was a visiting Ph.D student at University of Alberta, Edmonton, Canada. He is currently a Professor at the School of Traffic and Transportation Engineering, Central South University, Changsha, China.

His current research interests include Planning and Scheduling, Computational Intelligence and Machine Learning. He has authored more than 60 referred papers including those published in *IEEE TCYB*, *IEEE TSMCA*, *Information Sciences* and *Computers & Operations Research*. He serves as an Associate Editor of *Swarm and Evolutionary Computation Journal*, an editorial board member of *International Journal of Bio-Inspired Computation*, and a Guest Editor of *Information Sciences* and *Memetic Computing*. He is a regular reviewer of more than 20 journals including *IEEE TEVC*, *IEEE TCYB*, *Information Sciences* and *Applied Soft Computing*.



**Yingwu Chen** received the B.S. degree in automation, the M.S. degree in system engineering, and the Ph.D. degree in engineering from the National University of Defense Technology (NUDT), Changsha, China, in 1984, 1987, and 1994, respectively.

He was a Lecturer from 1989 to 1994, and an Associate Professor from 1994 to 1999 at NUDT. From 1999 to 2017, he was a Professor and the Director of the Department of Management Science and Engineering, College of Information Systems and Management, NUDT. Now, he has been a Distinguished Professor of the College of Systems Engineering, NUDT, where he focuses on management theory and its applications. He has authored more than 70 research publications. His current research interests include assistant decision-making systems for planning, decision-making systems for project evaluation, management decisions, and artificial intelligence.

Dr. Chen is the Editor of the Principle of System Engineering (Press of National University of Defense Technology), and Technology of Quantificational Analysis (China Renmin University Press).



**Witold Pedrycz** (IEEE Fellow, 1998) is Professor and Canada Research Chair (CRC) in Computational Intelligence in the Department of Electrical and Computer Engineering, University of Alberta, Edmonton, Canada. He is also with the Systems Research Institute of the Polish Academy of Sciences, Warsaw, Poland. In 2009 Dr. Pedrycz was elected a foreign member of the Polish Academy of Sciences. In 2012 he was elected a Fellow of the Royal Society of Canada. In 2007 he received a prestigious Norbert Wiener award from the IEEE

Systems, Man, and Cybernetics Society. He is a recipient of the IEEE Canada Computer Engineering Medal, a Cajastur Prize for Soft Computing from the European Centre for Soft Computing, a Killam Prize, and a Fuzzy Pioneer Award from the IEEE Computational Intelligence Society.

His main research directions involve Computational Intelligence, fuzzy modeling and Granular Computing, knowledge discovery and data science, pattern recognition, data science, knowledge-based neural networks, and control engineering. He has published numerous papers in these areas; the current h-index is 107 (Google Scholar). He is also an author of 18 research monographs and edited volumes covering various aspects of Computational Intelligence, data mining, and Software Engineering.

Dr. Pedrycz is vigorously involved in editorial activities. He is an Editor-in-Chief of *Information Sciences*, Editor-in-Chief of *WIREs Data Mining and Knowledge Discovery* (Wiley), and Co-editor-in-Chief of *Int. J. of Granular Computing* (Springer) and *J. of Data Information and Management* (Springer). He serves on an Advisory Board of *IEEE Transactions on Fuzzy Systems* and is a member of a number of editorial boards of international journals.