# Reinforcement Learning for Scalable Train Timetable Rescheduling with Graph Representation

Peng Yue, Yaochu Jin, *Fellow, IEEE,* Xuewu Dai, *Member, IEEE,* Zhenhua Feng, *Senior Member, IEEE,* Dongliang Cui

*Abstract*—Train timetable rescheduling (TTR) aims to promptly restore the original operation of trains after unexpected disturbances or disruptions. Currently, this work is still done manually by train dispatchers, which is challenging to maintain performance under various problem instances. To mitigate this issue, this study proposes a reinforcement learning-based approach to TTR, which makes the following contributions compared to existing work. First, we design a simple directed graph to represent the TTR problem, enabling the automatic extraction of informative states through graph neural networks. Second, we reformulate the construction process of TTR's solution, not only decoupling the decision model from the problem size but also ensuring the generated scheme's feasibility. Third, we design a learning curriculum for our model to handle the scenarios with different levels of delay. Finally, a simple local search method is proposed to assist the learned decision model, which can significantly improve solution quality with little additional computation cost, further enhancing the practical value of our method. Extensive experimental results demonstrate the effectiveness of our method. The learned decision model can achieve better performance for various problems with varying degrees of train delay and different scales when compared to handcrafted rules and state-of-the-art solvers.

*Index Terms*—Train timetable rescheduling, reinforcement learning, state representation, graph neural network

## I. INTRODUCTION

**H**IGH-speed railways provide a great convenience for passengers traveling medium and long distances [1]. However, many unexpected events may disturb or disrupt the normal operation of the trains, thus reducing the operation efficiency. For example, a minor malfunctioning of infrastructure tends to slightly defer the nominal arrivals and departures of multiple trains in the network, denoted as train disturbances. In extreme weather, long delays will occur and strongly affect the nominal circulation of the trains, which is defined as train disruptions [2]. To deal with these delay scenarios, human train dispatchers must modify the train schedules properly so that

Peng Yue and Dongliang Cui are with the State Key Laboratory of Synthetical Automation for Process Industries, Northeastern University, Shenyang 110819, China (e-mail: pyue515@gmail.com; cuidongliang@mail.neu.edu.cn).

Yaochu Jin is with the State Key Laboratory of Synthetical Automation for Process Industries, Northeastern University, Shenyang 110819, China. He is also with the Department of Computer Science, University of Surrey, Guildford, GU2 7XH, UK (e-mail: yaochu.jin@surrey.ac.uk) *(Corresponding author: Yaochu Jin)*

Xuewu Dai is with Northumbria University, Newcastle upon Tyne NE1 8ST, U.K., and also with Northeastern University, Shenyang 110819, China (e-mail: xuewu.dai@northumbria.ac.uk)

Zhenhua Feng is with the School of Computer Science and Electronic Engineering, University of Surrey, Guildford, GU2 7XH, UK. (E-mail: z.feng@surrey.ac.uk)

the trains can keep running in an efficient and ordered manner. This procedure is known as train timetable rescheduling (TTR) [3]–[5]. Currently, although computer-aided systems, *e.g.*, the centralized traffic control (CTC) system in China, can reduce the workload of dispatchers, such as train operation monitoring and command sending, it is still necessary for dispatchers to manually provide a new feasible schedule based on their experience [6]. Usually, it is extremely challenging for dispatchers to accomplish high-performance rescheduling when dealing with various delay scenarios. In this case, the design of an effective TTR algorithm is crucial not only for improving the efficiency of railway operations but also for reducing the burden of human dispatchers.

Like most scheduling problems, the TTR problem is NP-hard [7]. For this reason, a large body of research has been carried out to deal with this problem, which can be broadly classified into the following four categories. The first category is based on mathematical programming. Due to its theoretical guarantee of the solution quality, various methods have been proposed in this category [8]–[11]. Normally, they are pretty efficient when solving small-scale problems, but their solution time typically increases exponentially with the size of the problem, hindering their practicality to a certain degree. Many methods have been proposed to mitigate this issue, such as problem decomposition [12], [13], parallelization [5], and heuristics-assisted methods [14], [15]. The second approach to addressing this issue is to design meta-heuristic methods, such as genetic algorithms [16], tabu search [17] and hybrid algorithms [18]. Inspired by biological intelligence, these algorithms typically exhibit a strong global search capability, enabling them to find satisfactory solutions with a shorter response time compared to the first category of methods.

Despite the progress made in the two approaches discussed above, heuristic rules, which belong to the third category, remain prevalent in practice due to their higher computational efficiency. Popular heuristic rules include First-Come-First-Service (FCFS) [19] and First-Schedule-First-Service (FSFS) [20]. However, these heuristic rules heavily rely on domain knowledge in TTR, which leaves much room for improvement. Finally, learning-based methods, in particular reinforcement learning (RL), which are the fourth category, have been proposed for TTR in recent years. It has been proved that learning-based methods can learn advanced heuristic rules automatically, which not only significantly improves the solution quality but also maintain relatively low computational cost. Out of the above reasons, this work focuses on the RL-based approach to TTR. Notably, different from existing RL-

based methods, this work contains the following main new contributions:

- We employ the graph representation to describe the TTR problem. In the graph, the directed edge depicts the train operation constraints of TTR. With the help of a graph neural network (GNN), an informative state can be extracted automatically, thus improving the decision-making performance.
- We develop a tree search procedure to generate solutions for TTR, bringing about two main benefits. First, it facilitates the decoupling of the model's parameters from the problem's scale. Second, the track capacity constraints intrinsic to TTR can be dealt with easily by pruning infeasible branches in the constructed tree, thus ensuring the feasibility of the final scheme.
- We propose a learning curriculum for model training, comprised of unbalanced instance sampling and knowledge distillation. By employing this approach, the learned model exhibits better *scalability*, effectively adapting to problems with differing degrees of train delays.
- We propose a simple yet efficient local search method for the TTR problem. By refining the solutions generated by the learned decision model, the quality of the solutions can be significantly improved with only a small additional computational cost, further enhancing the scalability of the proposed method.

The rest of this paper is organized as follows. The related work is introduced in Section II. Then we describe the TTR problem and model it as a Markov decision process (MDP) in Section III. Section V elaborates our designed network architecture, while the learning algorithm and local search are discussed in Section VI. We evaluate the proposed method and report the experimental results in Section VII. Last, the conclusion is drawn in Section VIII.

## II. RELATED WORK

Inspired by the success of Alpha Go, some researchers have attempted to apply the RL approach to learn dispatching policies for TTR. One pioneering work was proposed by Šemrov *et al.* [21], proving that the learned policies are capable of outperforming the handcrafted rules. Li *et al.* provided a novel scheme based on multi-agent deep reinforcement learning [22], where the learned model can generate an optimal solution in the observed instances. Although the above methods have yielded excellent results, some challenges still remain when they are applied to real-world scenarios. One key obstacle pertains to the scalability of the approach. Specifically, we require the learned model to effectively address problems of different scales while maintaining high performance even under various levels of train delays. Notably, recent research has taken specific measures to mitigate this concern.

To ensure the model's scalability across various problem sizes, a typical approach involves manual extraction of key features from the railway system to create fixed-size states. For example, Khadilkar *et al.* defined their state by only using track information near the train to be dispatched [23]. A similar approach is also applied in [24], [25]. However, these manually designed states heavily rely on domain expertise and might miss vital factors affecting TTR, such as operational conflicts between trains. Ning *et al.* attempted automated state extraction using a convolutional neural network [26], but they trained their model on a single instance, making it challenging to generalize the extracted state to other instances. To resolve the above problem, we designed a directed graph to represent the TTR problem. This approach allows for the storage of planned and actual train operation details in nodes, while the train operation conflicts are represented by directed edges in the graph. Then a GNN model is employed to extract features from the graph, which can provide more relevant information compared to manually designed features and are independent of the problem size.

To ensure the model's scalability to various delay scenarios, the decision model must gain prior exposure to these scenarios during the training process. Wang *et al.* have made attempts in this regard, and eventually, their trained model can be applied to simple delay scenarios [25], where only a single train is delayed throughout the entire timetable. Subsequently, Peng *et al.* introduced an algorithm capable of handling more complex delay scenarios, allowing for minor delays on each train [27]. However, to the best of our knowledge, there are currently no RL-based methods tailored to address train disruptions involving prolonged train delays. Our experimental results show the difficulty in training a decision model to effectively handle these diverse scenarios simultaneously. To deal with this issue, we design a learning curriculum for model training, allowing knowledge transfer from small delay scenarios to large delay scenarios, thereby enhancing the training efficiency.

To enhance the model's scalability, additional search is frequently employed in the existing studies to seek solutions with improved performance. A common approach is to utilize a sampling strategy, constructing multiple solutions by selecting actions based on the probability distribution learned by the model [28]. Bello *et al.* introduced an effective search method known as active search, where a trained model is retrained to adjust the sampling probability [29]. Subsequently, Hottung *et al.* proposed an improved version to reduce the computational burden by limiting the model's weight that needs retraining [30]. Despite the success of these sampling methods in maintaining scalability, they are not suited for addressing the TTR problem, particularly in the context of high-speed railways. This is due to the requirement of real-time response of TTR within seconds, making it impractical to adequately implement sampling methods. Another widely used approach in RL-based methods is local search [31]. In comparison to the aforementioned methods, local search offers a notable advantage in computational efficiency, making it a more suitable choice for the TTR problem. However, to the best of our knowledge, there is no local search algorithm specifically designed for the TTR problem. Hence, we propose a local search method tailored for TTR, aiming to ensure the model's scalability across various delay scenarios.

TABLE I
ALL NOTATIONS USED IN THE TTR PROBLEM

| Type | Symbol | Definition |
|---|---|---|
| Decision variables | $a_{k,i}$ | Arrival time of train $k$ at station $i$ |
| | $d_{k,i}$ | Departure time of train $k$ at station $i$ |
| | $e_{k,i}$ | Delay penalty of train $k$ at station $i$ |
| | $y_{k_1,k_2,i}$ | If train $k_1$ departs earlier than train $k_2$ at station $i$, $y_{k_1,k_2,i} = 1$, otherwise, $y_{k_1,k_2,i} = 0$. |
| | $z_{k,i,p}$ | If train $k$ occupies track $p$ at station $i$, $z_{k,i,p} = 1$, otherwise, $z_{k,i,p} = 0$. |
| Parameters | $I$ | Number of stations |
| | $K$ | Number of trains |
| | $a_{k,i}^*$ | Original arrival time of train $k$ at station $i$ |
| | $d_{k,i}^*$ | Original departure time of train $k$ at station $i$ |
| | $rd_i$ | Minimum dwell time of trains at station $i$ |
| | $h$ | Minimum headway between trains |
| | $P_i$ | Track capacity at station $i$ |
| | $rt_{k,i}$ | Minimum running time of train $k$ between station $i$ and $i + 1$ (*i.e.*, section $i$) |
| | $M$ | A sufficiently large positive value |
| | $\lambda$ | Penalty factor for the early arrival of trains |
| | $e_{k,i}^*$ | Train delays that have occurred for train $k$ at station $i$ |

We set the value of $M$ to be 10 times greater than the maximum value observed in the planned train schedule when using the Gurobi solver to solve the TTR problem.

## III. PROBLEM DESCRIPTION

### A. Problem Description

A mathematical description of the TTR problem is given in this subsection. The relevant variables are defined in Table I.

The main purpose of TTR is to recover normal train operation as soon as possible, which minimized the following objective function:

$$J = \sum_{i=1}^{I} \sum_{k=1}^{K} (e_{k,i}) \tag{1}$$

where $e_{k,i}$ denotes the delay penalty of train $k$ at station $i$, and it is defined as:

$$e_{k,i} = \begin{cases} a_{k,i} - a_{k,i}^* & \text{if } a_{k,i} > a_{k,i}^* \\ \lambda * (a_{k,i}^* - a_{k,i}) & \text{else} \end{cases} \tag{2}$$

where $a_{k,i}$ is the arrival time of train $k$ at station $i$, $a_{k,i}^*$ is the planned one of train $k$ at station $i$, $\lambda$ represents the relative weighting of early arrivals relative to late arrivals of trains[1]. It is noteworthy that this defined non-linear objective function can be linearized as:

$$\begin{cases} e_{k,i} \geqslant a_{k,i} - a_{k,i}^* \\ e_{k,i} \geqslant \lambda * (a_{k,i}^* - a_{k,i}) \end{cases} \tag{3}$$

Due to the limitation of physical conditions, train operations should satisfy the following constraints.

[1] $\lambda$ is set to 0.3 in the experimental section, consistent with the parameter setting in [32].

*1) The minimum operation time constraint:* When train $k$ passes through one station or one section (between two stations), a certain operation time will surely be spent:

$$\begin{cases} a_{k,i+1} - d_{k,i} \geqslant rt_{k,i}, & \forall i = 1, \ldots, I-1; k = 1, \ldots, K; \\ d_{k,i} - a_{k,i} \geqslant rd_i, & \forall i = 1, \ldots, I; k = 1, \ldots, K; \end{cases} \tag{4}$$

where $d_{k,i}$ is the departure time of train $k$ at station $i$, $rt_{k,i}$ is minimum running time of train $k$ in section $i$, $rd_i$ is minimum dwell time of trains at station $i$.

*2) The planned and delay constraints:* The trains should not depart earlier than their planned time to ensure the normal passenger flow, and due to unexpected events, some trains have to be delayed, which can be formulated as:

$$\begin{cases} d_{k,i} - d_{k,i}^* \geqslant 0, \\ a_{k,i} - a_{k,i}^* \geqslant e_{k,i}^*, \end{cases} \forall i = 1, \ldots, I; k = 1, \ldots, K; \tag{5}$$

where $d_{k,i}^*$ is the original departure time of train $k$ at station $i$, and $e_{k,i}^*$ is the train delay of train $k$ at station $i$.

*3) The headway constraints:* For these trains running in the same section, overtaking is not allowed due to the limited tracks, while a time interval $h$ should be maintained for any two trains, as described by the following constraints:

$$\begin{cases} a_{k_2,i+1} - a_{k_1,i+1} \geqslant h - M(1 - y_{k_1,k_2,i}), \\ a_{k_1,i+1} - a_{k_2,i+1} \geqslant h - M * y_{k_1,k_2,i}, \\ d_{k_2,i} - d_{k_1,i} \geqslant h - M(1 - y_{k_1,k_2,i}), \\ d_{k_1,i} - d_{k_2,i} \geqslant h - M * y_{k_1,k_2,i}, \\ \forall i = 1, \ldots, I; k = 1, \ldots, K; \end{cases} \tag{6}$$

where $y_{k_1,k_2,i} = 1$ means that train $k_1$ will depart station $i$ earlier than train $k_2$.

*4) The track occupation constraints:* For any two trains (*e.g.*, train $k_1$ and train $k_2$), when they intend to occupy the same track $p$ at station $i$, a minimum time interval $h$ is required:

$$a_{k_2,i} - d_{k_1,i} \geqslant h - M * (3 - y_{k_1,k_2,i-1} - z_{k_1,i,p} - z_{k_2,i,p}), \\ \forall i = 2, \ldots, I; k = 1, \ldots, K; \tag{7}$$

where $z_{k_1,i,p} = 1$ ($z_{k_2,i,p} = 1$) denotes train $k_1$ ($k_2$) will occupy the track $p$ at station $i$.

*5) Natural constraints:* Based on the variable definition, the defined variables should satisfy the following:

$$y_{k_1,k_2,i} \in \{0,1\}, \ \forall i = 1, \ldots, I; k_1, k_2 = 1, \ldots, K;$$

$$\begin{cases} a_{k,i} \geqslant 0, \\ d_{k,i} \geqslant 0, \end{cases} \forall i = 1, \ldots, I; k = 1, \ldots, K;$$

$$\begin{cases} z_{k,i,p} \in \{0,1\}, \\ \sum_{p=1,\ldots,P_i} z_{k,i,p} = 1. \end{cases}$$

$$\forall p = 1, \ldots, P_i; i = 2, \ldots, I; k = 1, \ldots, K; \tag{8}$$
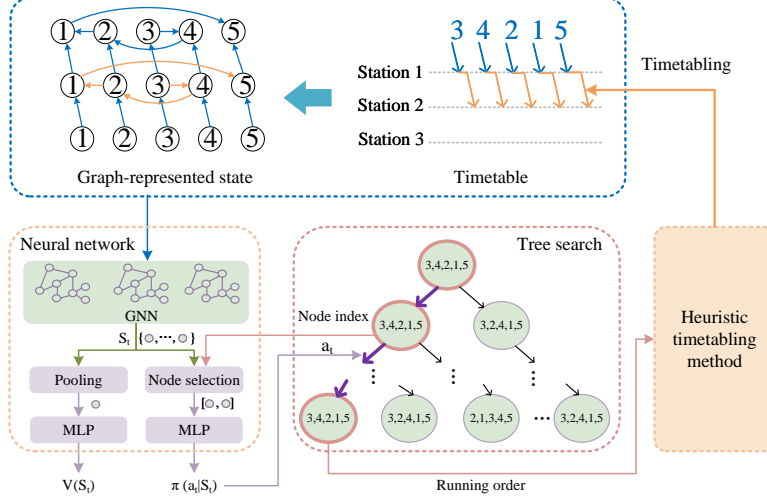
Fig. 1.   The designed MDP for constructing the TTR's solution.

## IV. MARKOV DECISION PROCESS FOR TTR

When applying the RL algorithm to solve the TTR problem, the first step is to build a suitable Markov decision process (MDP). In this study, we followed the approach outlined in [27], [33], where the decision model only needs to determine the running order between the trains, and the specific time is obtained by using existing heuristics. Figure 1 provides a detailed illustration of our designed decision process.

Specifically, a TTR's solution is constructed section by section. For each section, a timetable is first transformed into a directed graph. The GNN module then extracts a state $S_t$ from this graph. Subsequently, the tree search process is implemented, guided by the neural network, to determine the running order of trains. At each decision branch, action $a_t$ represents the routing selection. Finally, an existing heuristic method generates the operation scheme of the trains based on the determined running order. This iterative process continues until the whole timetable is constructed. Note that for clarity, we only show the detailed information needed at each stage in this figure, omitting the state transition and reward calculation during the tree search, which will be elaborated in Subsection V-C.

One thing we need to mention is that this designed decision process ensures the model's scalability in terms of the problem size. Firstly, the GNN module is inherently independent of the number of nodes in the graph (corresponding to the problem size). Then, during the tree search process, the fixed-length node embedding is obtained through node selection and fed into the Multi-Layer Perceptron (MLP) layer. This way, the parameters of the neural network are totally independent of the problem size.

Regarding the heuristic timetabling method, it is also employed in the past studies [27], [33], but we make some revisions to allow for early arrival of the trains. The details can be found in Algorithm 1.

In Algorithm 1, the train operation scheme for each train is determined one by one according to the given running order

---

**Algorithm 1** A heuristic timetabling method

**Input:**
   $O_a^i$: The given order of trains in section $i$
   $\{a_{k,i}$ for $k = 1, \ldots, K\}$: The arrival time of trains at station $i$

**Output:**
   $\{a_{k,i}, a_{k,i+1}$ for $k = 1, \ldots, K\}$: The updated arrival time
   $\{d_{k,i}$ for $k = 1, \ldots, K\}$ : The departure time of trains at station $i$

1: Sort the trains based on $O_a^i$ to obtain $\Omega_s$
2: **for** $k$ in $\Omega_s$ **do**
3:     Calculate the number of trains at station $i$ when train $k$ arrives at time $a_{k,i}$, noted as $n_k$
4:     **if** $n_k = P_i$ **then**
5:         Delay the arrive time $a_{k,i}$ of train $k$ by:
            $a_{k,i} = h + d_{g_1,i}$ where there are $(P_i - 1)$ trains between train $k$ and train $g_1$.
6:         Allocate the track occupied by train $g_1$ to train $k$.
7:     **else**
8:         Allocate the unoccupied track randomly to train $k$.
9:     **end if**
10:     Calculate the time $d_{k,i}$ and $a_{k,i+1}$ by Eq. (9) and (10).
11: **end for**
12: **for** $k$ in $\Omega_s$ **do**
13:     Calculate buffer time $B_h(g_2, k, i+1)$ by Eq. (11)
14:     **if** $B_h(g_2, k, i+1) < 0$ **then**
15:         Get train set $\Omega_p$ that can advance their arrival time.
16:         Advance arrival time of the train in $\Omega_p$ by Eq. (13).
17:         **if** $(\Delta_{g_2}) < |B_h(g_2, k, i+1)|$ **then**
18:             Delay arrival time $a_{k_f,i+1}$ by Eq. (16)
19:         **end if**
20:     **end if**
21: **end for**

$O_a^i$. Specifically, for train $k$, we first check whether there are available tracks for it to stop in station $i$. If not, its arrival time must be delayed until a train (corresponding to train $g_1$) leaves the station. It also meant that train $k$ and $g_1$ have to occupy the same track on station $i$, referring to lines 3-9 in Algorithm 1. For its departure time $d_{k,i}$, we employ equation (9) to obtain a feasible and earliest time, hoping to reduce its delay as much as possible.

$$d_{k,i} = \max\left(d_{k,i}^*, a_{k,i} + rd_i, d_{g_2,i} + h\right) \quad (9)$$

where train $g_2$ is the preceding train of train $k$. The same manner can also be used to obtain the arrival time at the next station [27], but it requires that the train can not arrive at stations earlier than its planned time, which is not necessary in practice. To fill this gap, we provide a new manner to obtain a more flexible solution. Specifically, we first calculate the earliest time for each train to arrive at the next station, ignoring the operating conflicts between trains, referring to Eq. (10).

$$a_{k,i+1} = \max(a_{k,i}^*, d_{k,i} + rt_{k,i}) \quad (10)$$

To validate whether the obtained time is feasible, we will calculate the buffer time $B_h(g_2, k, i+1)$ between trains[2].

$$B_h(g_2, k, i+1) = a_{k,i+1} - a_{g_2,i} - h \quad (11)$$

where train $g_2$ is the preceding train of train $k$. If there is a conflict between adjacent trains (i.e., $B_h(g_2, k, i + 1) < 0$), minor revisions will be implemented to dissipate conflicts, referring to lines 15-19.

There are two main steps to complete the revisions. The first step is to advance the arrival time of the trains in $\Omega_p$, denoting a set of trains before a conflict occurs[3]. For this reason, we need to calculate the adjustable time for these trains iteratively while satisfying their own constraints, calculated by Eq. (12).

$$h_{aj}(k_p) = \min\left(B_r(k_p, i), h_{aj}(k_p\text{-}1) + B_h(k_p\text{-}1, k_p, i+1)\right) \quad (12)$$

where $B_r(k_p, i) = a_{k_p,i+1} - a_{k_p,i} - rt_{k,i}$ represents the supplement time of train $k_p$ in segment $i$. For the earliest arrival train in $\Omega_p$, $h_{aj}(k_p)$ is set as $B_r(k_p, i)$. Next, we advance the arrival time of these trains in $\Omega_p$ by $\Delta_{k_p}$:

$$a_{k_p,i+1} = a_{k_p,i+1} - \Delta_{k_p} \quad (13)$$

Specifically, for train $g_2$, i.e., the preceding train involved in conflicts, its actual adjustment $\Delta_{g_2}$ is the minimum value between the adjustable time $h_{aj}(g_2)$ and the violation $(-B_h(g_2, k, i + 1))$:

$$\Delta_{g_2} = \min\left(-B_h(g_2, k, i+1), h_{aj}(g_2)\right) \quad (14)$$

To avoid triggering new conflicts, we also need to adjust other preceding trains in $\Omega_p$ iteratively:

$$\Delta_{k_p\text{-}1} = \max\left(\Delta_{k_p} - B_h(k_p\text{-}1, k_p, i+1), 0\right) \quad (15)$$

[2]Here, we only need to use the buffer time to check if the minimum interval is satisfied between trains (i.e., headway constraints), since the constraints involving a single train have already been considered in Eq. (10).

[3]The number of the trains in $\Omega_p$ is determined by $\lambda$. For example, when $\lambda$ is 0.3, then the number of trains that can be advanced is 3, as it incurs a smaller cost in objective $J$ to delay one train than to advance four trains.
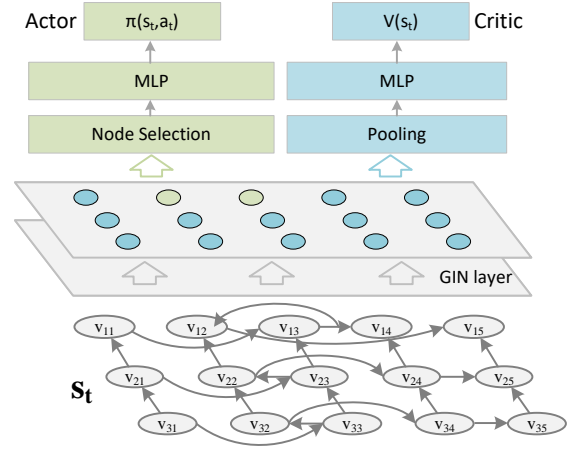


Fig. 2. The network architecture for the decision model.

If the above adjustment is not enough to resolve the conflict (i.e., $\Delta_{g_2} < |B_h(g_2, k, i+1)|$), the second step will be implemented to delay the arrival time of the trains in $\Omega_s(k\text{:end})$ iteratively, denoting the set of trains after the conflict occurs. The specific updating process is shown in Eq. (16).

$$a_{k_f,i+1} = \max\left(a_{k_f,i+1}, (a_{k_f\text{-}1,i+1} + h)\right) \quad (16)$$

where train $k_f \in \Omega_s(k\text{:end})$ indicates the train operating after the conflict between $g_2$ and $k$.

## V. NETWORK ARCHITECTURE

In this section, we will initially present the designed network architecture for TTR. Subsequently, we will provide a detailed description of its application in constructing solutions for TTR, encompassing the GNN-based state and tree search. Additionally, we will elaborate on the state transitions and how rewards are calculated during the tree search process.

The whole network architecture is shown in Fig. 2, which can be divided into two parts: the GIN layer and the task-specific layer. Firstly, the GIN layer is applied to extract the state features. Subsequently, the critic subnetwork undertakes the state evaluation, denoted as $V(S_t)$, while the actor subnetwork is engaged in action selection $\pi(a_t|s_t)$. These components will be thoroughly elucidated in the subsequent subsection.

### A. GNN-Based State Representation

In this subsection, we will describe how to describe a TTR instance as a directed graph, thus extracting the state information by using the GIN layer.

*1) Graph Representation for TTR:* In the TTR problem, it is challenging for the dispatcher to determine an effective running order between trains. The main reason is that there are so many factors influencing the scheduling results, *e.g.*, the original schedule, the actual schedule, the type of locomotive, and the coupling between trains, so it is also non-trivial
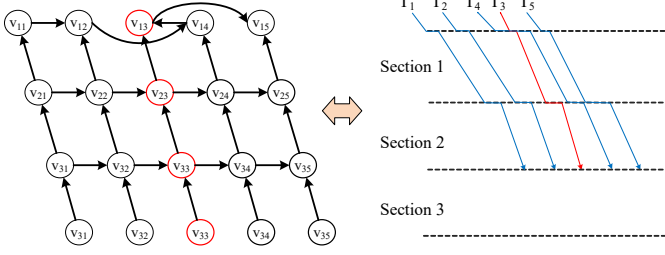
Fig. 3. An example of modeling the TTR as a directed graph.



Fig. 4. An example of an infeasible order.

for a handcrafted state to grasp this information completely. Therefore, we attempt to describe this information by using a directed graph, and a GNN can be used to grasp key information automatically. Figure 3 shows an example of a designed graph.

The graph can be defined as $G = (V, E)$, which consists of two components: nodes $V = \{v_{k,i} : [\delta_{k,i}, a^*_{k,i}, L_{k,i}]\}$ and edges $E$. Node $v_{k,i}$ represents the arrival event of train $k$ at station $i$ and records three features, including train delay $\delta_{k,i}$ (calculated by Eq. (17)), planned arrival time $a^*_{k,i}$ and a binary indicator $L_{k,i}$ indicating whether it has been rescheduled ($L_{k,i} = 1$ means it has been rescheduled).

$$\delta_{k,i} = \begin{cases} \max\big(0, \delta_{k,i-1} - B_r(k, i\text{-}1)\big) & \text{if } L_{k,i} = 0 \\ \text{Calculated by Algorithm 1} & \text{if } L_{k,i} = 1, \end{cases} \quad (17)$$

where $B_r(k, i\text{-}1)$ indicates the supplement time of train $k$ in section $(i\text{-}1)$. It can be found that the train delay is not always accurate when $L_{k,i}$ is 0, but it is efficient in improving the training efficiency.

The edge connecting two nodes describes the constraint relationship between adjacent events, which is similar to the existing AG graph. But there exist significant differences between them. In the AG graph, the edge direction is consistent with the sequence of events, proceeding from the antecedent event to the subsequent one. Conversely, in our designed graph, when two adjacent events describe the same train at different stations, we reverse the edge's direction, making it point from the subsequent event to the antecedent one (*e.g.*, the edge between node $v_{13}$ and $v_{23}$). This modification enables the decision model to obtain future information from subsequent nodes, thus better determining the running order of trains.

*2) State Extraction:* Based on the graph, a GNN model can be used to extract the state feature automatically. In the GNN family, we adopt the Graph Isomorphism Network (GIN) [34] to achieve this due to its strong discriminative power.

Specifically, given a graph $G = (V, E)$, the node representation can be extracted in an iterative way. In each iteration, the node embedding is updated by aggregating the information from its neighborhoods, which is formulated as:

$$h_v^{(k)} = \text{ReLU}\left(\text{BN}\left(\text{MLP}_{\theta_k}^{(k)}\left((1+\epsilon^{(k)})h_v^{(k-1)} + \sum_{u \in \mathcal{N}(v)} h_u^{(k-1)}\right)\right)\right), \quad (18)$$

where $h_v^{(k-1)}$ indicates the node embedding of node $v$ in the iteration $(k\text{-}1)$, and the raw feature $h_v^0$ is defined as
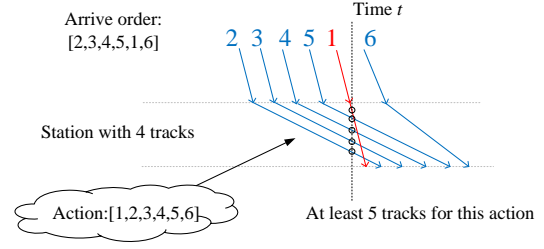
$h_v^0 = [\delta_{k,i}, a^*_{k,i}]$. $\mathcal{N}(v)$ is the neighborhoods of the nodes $v$, $\text{MLP}_{\theta_k}^{(k)}(\cdot)$ is a multi-layer perceptron (MLP) with parameter $\theta_k$, $\epsilon^{(k)}$ is learnable scalar, $\text{BN}(\cdot)$ stands for batch normalization, $\text{ReLU}(\cdot)$ is the rectified linear unit.

### B. Running order determination by using tree search process

Given the extracted node embedding, we need to provide a feasible running order for trains. A prevalent method for sequence determination in prior studies [33], [35] involves incrementally incorporating one item (i.e., trains) into the partial sequence, but this process ignores the inherent constraints specific to TTR, denoted as track capacity constraints. Figure 4 shows a simple example of an infeasible running order.

It can be found that each train is not allowed to overtake more than $(P_i - 1)$ trains at station $i$, otherwise, the number of trains will exceed the number reserved. To address this issue, this study designs a tree search process to obtain feasible solutions for TTR, as shown in Fig. 5.

The running order between trains in a section can be determined by getting a route in the constructed tree, and the task of the decision model is thus transformed into giving priority between two adjacent trains at each bifurcation point. There are mainly two benefits. First, the tracking capacity constraints can be addressed directly by pruning the infeasible branches. Second, the fixed-length node embedding is obtained by concatenating the two nodes related to the trains to be dispatched, thus decoupling the model's parameters from the problem size. Since this procedure is similar to the dispatcher's
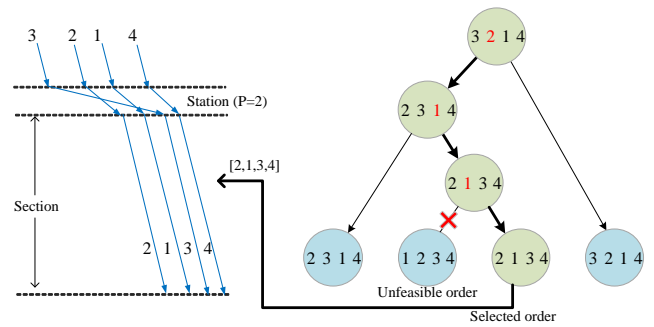


Fig. 5. The tree search process for generating a feasible running order of trains.

**Algorithm 2** Construction of the search tree

**Input:**
    $O_a^{i-1}$: The actual running order of trains in section $i$-1;
    $O_p^i$: The planned running order of trains in section $i$

**Output:**
    The constructed tree;
1:  $O_0 = O_a^{i-1}$       % The root node of the tree
2:  **% Determine the overtaken trains** $\Omega_{ot}$
3:  $O_{ac} = \text{copy}(O_a^{i-1})$
4:  **for** $k$ in $O_p^i$ **do**
5:     Find the location of train $k$ in $O_{ac}$:
        $loc_k = \text{find}(O_{ac}, k)$
6:     **if** $loc \neq 1$ **then**
7:         Add $k$ into set $\Omega_{ot}$
8:     **end if**
9:     Delete $k$ from $O_{ac}$
10:  **end for**
11:  **% Create the child nodes iteratively**
12:  $O_f = O_0$     % $O_f$ is the father node
13:  **for** $k$ in $O_a^{i-1}$ **do**
14:     **if** $k$ in $\Omega_{ot}$ **then**
15:         Find the proceeding train of train $k$ in $O_f$:
            $g_2 = O_f(\text{find}(O_f, k) - 1)$
16:         **while** $g_2$ is not empty **do**
17:             **if** $\text{find}(O_p^i, g_2) < \text{find}(O_p^i, k)$ **then**
18:                 **Break**
19:             **else**
20:                 Create two child nodes $O_{c1}$ and $O_{c2}$:
                    $O_{c1} = O_f$
                    $O_{c2} = \text{Swap}(O_f, k, g_2)$
21:                 $O_f = O_{c2}$
22:                 $g_2 = O_f(\text{find}(O_f, k) - 1)$
23:             **end if**
24:         **end while**
25:     **end if**
26:  **end for**

mitigates the effects on passengers who require transfers.

Recall that the decision model's task is to select a route at each decision branch. Thus, we design the upper layer of the neural network as follows:

For the actor subnetwork, the concatenation operator $[\cdot, \cdot]$ is firstly implemented in the node selection layer for two nodes $v_{i,k_1}$ and $v_{i,k_2}$ to be compared, and an MLP layer and sigmoid function $\sigma(\cdot)$ are then applied to determine their priorities:

$$\pi_\psi(a_t|s_t) = \sigma(\text{MLP}_\psi([h_{v_{i,k_1}}^K, h_{v_{i,k_2}}^K])). \tag{19}$$

where $h_{v_{i,k_1}}^K$ and $h_{v_{i,k_2}}^K$ are the selected node embedding.

For the critic subnetwork, a global representation $h_g$ of the entire graph is obtained by using the average pooling layer. We apply an MLP to obtain the state value $V(s_t)$:

$$\begin{cases} h_g = \text{Pool}(\{h_v^K : v \in V\}) = \sum_{v \in V} h_v^K / |V| \\ V(s_t)_\phi = \text{MLP}_\phi(h_g). \end{cases} \tag{20}$$

### C. State Transition and Reward Definition

*1) State Transition:* At each decision step, the decision model will provide a priority between two adjacent trains. If the given order is changed, the constraint relationships are also updated, as reflected by the change of edge in the graph. Once the decision model has completed a route search of the built tree in a section, an entire running order of the trains can be obtained. At this point, it is necessary to update the node features to provide a more precise prediction of train delays (*i.e.*, $\delta_{k,i}$), using Algorithm 1.

For easier understanding, Fig. 6 depicts the process of state transitions. At the decision step $t$, the decision model provides the priority between trains $T_2$ and $T_4$, and the final decision is

behavior, we employ the dispatcher's knowledge to construct the tree, thereby enhancing the search efficiency. The specific process is illustrated in Algorithm 2.

First, we define the running order of trains in the last section (*i.e.*, $O_a^{i-1}$) as the root node $O_0$ in the tree. Second, the overtaken trains $\Omega_{ot}$ are identified iteratively by comparing the positions of the same trains in $O_a^{i-1}$ and $O_p^i$, referring to lines 4-10. Third, we create the child nodes in the tree iteratively by swapping the position between overtaken train $k$ and its preceding train $g_2$ (*i.e.,* $\text{Swap}(O_f, k, g_2)$). It should be noted that train $k$ cannot overtake the train with a higher priority in planned order $O_p^i$, referring to lines 17-18. When all overtaken trains in $\Omega_{ot}$ are adjusted, the whole tree is constructed. It can be found that the constructed tree will help the overtaken trains to restore their original positions as soon as possible. The motivation stems from the actions of train dispatchers. In practice, train dispatchers have a preference for selecting the departure order closer to the original one, which is due to the fact that it has less impact on the subsequent scheduling tasks, such as the rolling stock scheduling [36]. Furthermore, it also
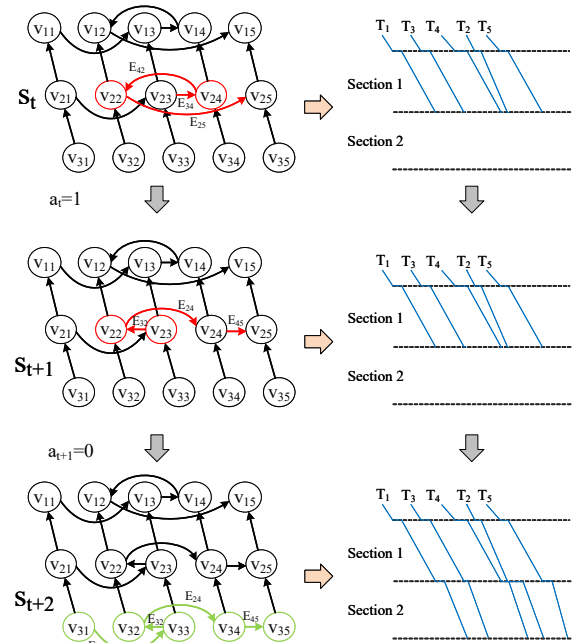


Fig. 6. An example of state transition.

to reverse their order (*i.e.*, $a_t = 1$), so related constraints have to be changed. For example, the red edges $E_{25}$ (from node $v_{22}$ to $v_{25}$), $E_{42}$ and $E_{34}$ in the graph are deleted, and some new edges are added, including $E_{32}$, $E_{24}$ and $E_{45}$. In the next step, the decision model chooses to maintain the priority between trains $T_2$ and $T_3$ (*i.e.*, $a_{t+1} = 0$), so the corresponding edge remains unchanged. Since the entire tree search procedure has been completed, Algorithm 1 can be applied to update the train schedule in this section. Accordingly, some new edges (*i.e.*, green edges) reflecting the running order will be added to the graph, and the node features also need to be updated. To be more specific, $L_{w,i}$ for the green nodes is assigned a value of 1, and $\delta_{k,i}$ is calculated using Eq. (17).

*2) Reward Definition:* In this study, we define an unbiased reward for objective $J$, which is calculated by Eq. (21).

$$
r_t = \begin{cases} (\hat{J}^{t-1} - \hat{J}^t)/(K * I) & \text{if node feature is updated} \\ 0 & \text{else,} \end{cases}
$$

(21)

where $\hat{J}^{t-1}$ and $\hat{J}^t$ are the predicted objective value in the current state and the updated state, respectively, and $K * I$, as the description of problem size, is used to reduce the variance of reward. When the entire running order is determined by the decision model, the reward is defined as the decrease in the objective function $J$. Otherwise, the reward is assigned a value of 0.

## VI. LEARNING ALGORITHM AND LOCAL SEARCH

### A. Learning Algorithm

In this study, we aim to train a generic model scalable to various delay scenarios, including train disturbances and disruptions. However, our experiments reveal that current RL algorithms fall short in training such a scalable decision model. To address this, we design a learning curriculum (LC) to improve the training performance. The detailed implementation can be found in Algorithm 3.

We divide the training process into two stages. In the first stage, we employ the vanilla PPO algorithm to train the decision model on small delay scenarios. Then, in the second stage, the learned model undergoes retraining using our proposed learning curriculum to address the large delay scenarios. Our learning curriculum improves upon the vanilla PPO method in two key aspects: instance sampling and loss definition. The rationale behind this is to transfer knowledge obtained from handling small delay scenarios (a relatively easy task) to effectively handling large delay scenarios (a more challenging task).

Regarding instance sampling, the PPO algorithm generates TTR instances randomly, assigning each train a uniform random delay. In contrast, we introduce an unbalanced sampling weight $P_{LC} = [p_s, p_l](p_s > p_l)$, where $p_s$ is the sampling probability for small delay scenarios (train disturbances), and $p_l$ corresponds to large delay scenarios (train disruptions), as outlined in lines 3-7 of Algorithm 3. Concerning the loss definition, we introduce an additional loss term $L_4(\theta, \phi)$ to implement knowledge distillation, referring to lines 19-22 of

Algorithm 3. Specifically, the original loss function $L(\theta, \psi, \phi)$ used in PPO algorithm is defined as:

$$
\begin{cases} L_1(\theta, \psi) = -\sum_t^T \min\left(\mathcal{E}_t \hat{A}_t, clip(\mathcal{E}_t, 1 - \epsilon, 1 + \epsilon)\hat{A}_t\right) \\ L_2(\theta, \phi) = \frac{1}{T}\sum_t^T (\sum_t^T \gamma^t r_t - V(s_t))^2 \\ L_3(\theta, \psi) = -\sum_t^T H([\pi_{\theta,\psi}(s_t), 1 - \pi_{\theta,\psi}(s_t)]) \\ L(\theta, \psi, \phi) = w_1 L_1(\theta, \psi) + w_2 L_2(\theta, \phi) + w_3 L_3(\theta, \psi), \end{cases}
$$

(22)

where $\mathcal{E}_t = \frac{\pi(a_t|s_t)}{\pi_{old}(a_t|s_t)}$ is the probability ratio between the old policy $\pi_{old}$ and current policy $\pi$, and $\hat{A}_t = \sum_t^T \gamma^t r(s_t, a_t) - V(s_t)$ is an estimator of the advantage function at timestep $t$, and $H(\cdot)$ is the entropy loss function. In this function, $L_1(\theta, \psi)$ is responsible for improving the policy steadily, and $L_2(\theta, \phi)$ is used to estimate the state value $V(s_t)$, and $L_3(\theta, \psi)$ is employed to ensure the exploratory power of the model.

The extra loss $L_4(\theta, \psi)$ is defined as follows:

$$
L_4(\theta, \psi) = H(\pi_S(a_t|s_t), \pi(a_t|s_t)) \tag{23}
$$

where $\pi_S(a_t|s_t)$ represents the learned policy from the first stage, and $\pi(a_t|s_t)$ represents the current policy.

### B. Local Search for TTR

Using the learning curriculum and tree search, the learned decision model can address any TTR instances directly logically, involving various degrees of train delays and different

---

**Algorithm 3** Proposed Learning algorithm

**Input:** $\{\theta, \psi, \phi\}$: The initial parameter of the decision model
**Output:** Trained decision model $\{\theta^*, \psi^*, \phi^*\}$
1: **for** $stage\ sta = 1, \cdots, 2$ **do**
2:     **for** $episode\ epi = 1, \cdots, \text{Max\_epi}$ **do**
3:         **if** $sta == 1$ **then**
4:             Generate a problem instance randomly
5:         **else**
6:             Generate a instance based on distribution $P_{LC}$
7:         **end if**
8:         **while** True **do**
9:             Transform the TTR instance as a graph $G$
10:             State extraction by using GIN layers
11:             Get a running order $a_t$ by using the tree search
12:             Generate a partial solution using Algorithm 1
13:             **if** the whole timetable is generated **then**
14:                 **break**
15:             **end if**
16:         **end while**
17:         **for** $Epoch\ epo = 1, \cdots, \text{Max\_epo}$ **do**
18:             Calculate loss item $L(\theta, \psi, \phi)$ using Eq (22).
19:             **if** $sta == 2$ **then**
20:                 Calculate loss item $L_4(\theta, \phi)$ using Eq (23).
21:                 $L(\theta, \psi, \phi) = L(\theta, \psi, \phi) + w_4 \cdot L_4(\theta, \phi)$
22:             **end if**
23:             Update $\{\theta, \psi, \phi\}$ by using Adam optimizer
24:         **end for**
25:     **end for**
26: **end for**

problem sizes. However, maintaining peak performance across all instances remains challenging. To enhance the model's performance, we introduce a local search method during the application stage. This method is akin to the commonly used 2-opt heuristics [37] in TSP, as detailed in Algorithm 4.

The local search is mainly divided into three steps. First, we randomly select an arrival event of the train from the entire timetable, denoted as $v_{i,k}$. If train $k$ swaps its positions with its adjacent train $adj(k)$ compared to the previous station $i-1$ (i.e., $y_{k,adj(k),i} \neq y_{k,adj(k),i-1}$), the second step will be implemented to reorder train $k$ and $adj(k)$, making it consistent with the order at station $i-1$. Subsequently, Algorithm 1 is employed again to revise the specific arrival/departure time due to the reordering operation. Notably, this revision procedure only starts after the adjustment event and stops when the conflict is resolved, thus saving the computational cost. The change in objective function $J$ is then calculated and denoted as $\Delta J$, corresponding to lines 8-9 in Algorithm 4. The final step is to do further exploration if the performance of the scheme is improved, i.e., $\Delta J < 0$. To be more specific, the same reordering operation will be performed again between train $k$ and $adj(k)$ at subsequent stations until the solution quality cannot be further improved, see lines 10-19 of the Algorithm. To help the reader understand its procedure, we also give a simple example in Fig. 7.

In this example, node $v_{23}$ is firstly selected by executing the first step, since its position has changed compared to the first station. Next, we proceed to the second step to rearrange its order with its adjacent train $v_{24}$, depicted by the red edge in the
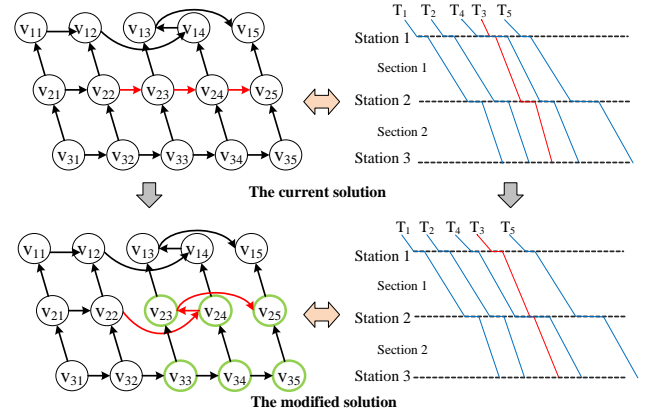


Fig. 7. An example of local search.

graph. By using Algorithm 1, the arrival time can be revised, denoted by the green nodes of the graph. If the revised scheme exhibits improved performance, we proceed to exchange the order between nodes $v_{33}$ and $v_{34}$ in the third step. If not, one local search iteration is completed.

## VII. EXPERIMENTAL RESULTS

### A. Datasets and Experimental Settings

To make it practical for real-time decision-making, the decision model must be pre-trained on diverse problem instances, enabling it to address unexpected disturbances rapidly. However, due to the privacy concerns associated with railway operation data, obtaining an adequate and representative set of delay scenarios is challenging. In this study, we generate augmented instances from limited instances by using Algorithm 5.

Specifically, we first construct a raw instance $TS_b$ by introducing random noise into train timetable $TS_a$. Then we design a matrix $M_{eq} = J_I \cdot \text{diag}(\text{range}(max(TS_b)))$, where $J_I$ is all-one matrix, and range($max(TS_b)$) stands for the arithmetic sequence with an upper bound as $max(TS_b)$. By adding this matrix into $TS_b$, the trains will be uniformly distributed over schedule $TS_b$, which is a common feature of the train

---

**Algorithm 4** Local search

**Input:**
    The search times Max_ls
    A rescheduled timetable.
**Output:**
    A updated rescheduled timetable
1: **for** $ls = 1, \cdots, \text{Max\_ls}$ **do**
2:     **while** True **do**
3:         Randomly select a overtaken train $k$ and station $i$.
4:         **if** $y_{k,adj(k),i} \neq y_{k,adj(k),i-1}$ **then**
5:             **Break**
6:         **end if**
7:     **end while**
8:     Order swap between train $k$ and $adj(k)$ at station $i$.
9:     Implement Algorithm 1 locally to rapidly estimate the objective change, denoted as $\Delta J$.
10:     **if** $\Delta J < 0$ **then**
11:         **for** $i_2 = i, \ldots, I$ **do**
12:             Order swap between train $k$ and $adj(k)$ at station $i_2$.
13:             Estimate the new objective change $\Delta J$.
14:             **if** $\Delta J > 0$ **then**
15:                 **Break**
16:             **end if**
17:         **end for**
18:     **end if**
19: **end for**

---

**Algorithm 5** Data augmentation for TTR instances

**Input:**
    An actual train schedule $TS_a$ for TTR;
    Disturbance parameters $[\tau_1, \tau_2, \tau_3]$
**Output:**
    The problem instances $TS_b$
1: $TS_b = TS_a + M_r * \tau_1$
2: $M_{eq} = J_I \cdot \text{diag}(\text{range}(max(TS_b)))$
3: $TS_b = TS_b + M_{eq}$
4: Calculate the running order $O_b$ based on $TS_b$
5: $TS_b$=Algorithm 1($TS_b, O_b$)
6: Assign an initial delay $ad \in [0, \tau_2]$ to each train, and set the minimum operation time $r_{k,t} \in [\tau_3 \cdot r_t^*, r_t^*]$
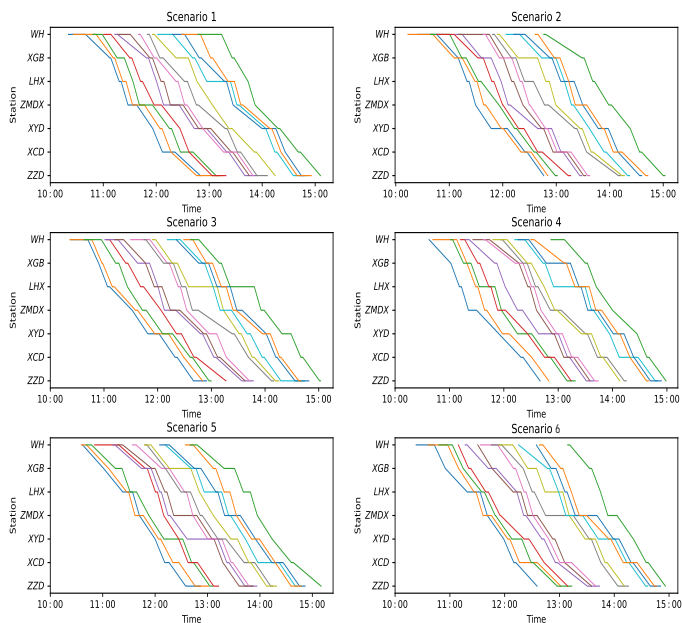
Fig. 8. The augmented TTR instances.

schedule. To obtain a feasible train schedule, Algorithm 1 is adopted to eliminate the conflicts in $TS_b$. Regarding the train delay, we assign a random delay $ad \in [0, \tau_2]$ to each coming train, and a minimum running time $r_{k,t} \in [\tau_3 \cdot r_t^*, r_t^*]$ is also set randomly to reflect the differences in train types, where $r_t^*$ is the original parameters in $TS_a$.

According to the real scenarios, the disturbance parameter $\tau_1$ is set as 20 (min), which can ensure the diversity of the original schedule as much as possible, and some of the generated examples are shown in Fig. 8. $\tau_2$ is set to 60 and 180 for train disturbances and train disruptions, respectively. Parameter $\tau_3$ is set to 0.3.

The hyperparameters of the proposed algorithm are set as presented in Table II. Notably, we assign parameter $w_3$ different values for different scenarios. Given that the values of $L_1$ and $L_2$ increase with the severity of train delays, a larger value of $w_3$ is required to balance the exploration in large-delay scenarios. We use a computer to conduct all experiments with a specification of an Intel Core i9-10940X CPU and 32GB RAM. The method is implemented in Python with PyTorch.

Recall that the model training process is divided into two stages. In the first stage, the vanilla PPO algorithm is employed to train the decision model under train disturbances. Subsequently, in the second stage, the learned model undergoes retraining using our proposed learning curriculum to handle train disruptions. Therefore, in the subsequent section, we will verify the algorithm's effectiveness step by step in this order.

## B. TTR for Train Disturbances

Given that the problem size (i.e., the number of decision variables) mainly depends on the number of stations ($I$) and trains ($K$), we denote the problem size as $I * K$. To verify the scalability of our approach, we first train our model on

TABLE II
THE HYPERPARAMETERS USED IN THE PROPOSED METHOD

| Type | Hyperparameter | Value |
|---|---|---|
| Network | Number of GIN layers | 1 |
| | Hidden dimensions of GIN layers | 128 |
| | Number of MLP layers for Critic | 2 |
| | Hidden dimension of MLP for Critic | 128 |
| | Number of MLP layers for Actor | 2 |
| | Hidden dimension of MLP for Actor | 128 |
| Training | Discount rate $\gamma$ | 0.9 |
| | Clipping rate $\epsilon$ | 0.2 |
| | The weights $w_1, w_2, w_3, w_4$ | [2, 2, (0.1, 0.03), 100] |
| | Number of training episodes | 5e4 |
| | Learning rate | 1e-4 |
| | Optimizer | Adam |
| | Number of updates for one episode | 10 |
| | Unbalanced sampling weight | [0.8, 0.2] |
| Application | Number of local search | 50 |



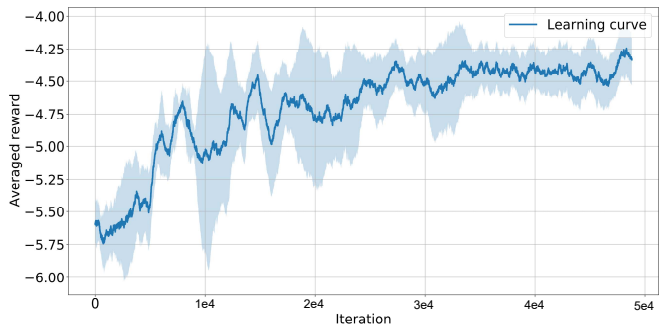Fig. 9. The learning curve of our model under train disturbances.

TABLE III
THE TTR RESULTS UNDER TRAIN DISTURBANCES

| Methods | Metrics | Test 10*10 | Generalization 15*15 | 20*20 | 20*30 |
|---|---|---|---|---|---|
| FCFS | Obj.(min) | 1064 | 2295 | 3262 | 5756 |
| | Gap | 68.42% | 76.30% | 80.41% | 88.67% |
| | Time(s) | 0.086 | 0.165 | 0.270 | 0.395 |
| FSFS | Obj.(min) | 657 | 1212 | 1361 | 1817 |
| | Gap | 48.86% | 55.12% | 53.05% | 64.12% |
| | Time(s) | **0.071** | **0.117** | **0.133** | **0.249** |
| RLG | Obj.(min) | 403 | 670 | 807 | 842 |
| | Gap | 16.63% | 18.81% | 20.82% | 22.57% |
| | Time(s) | 0.45 | 1.598 | 2.044 | 3.379 |
| RLG-LC | Obj.(min) | 374 | 610 | 753 | 769 |
| | Gap | 10.16% | 10.82% | 15.14% | 15.21% |
| | Time(s) | 2.176 | 3.131 | 3.737 | 4.295 |
| Gurobi (Baselines) | Obj.(min) | **336** | **544** | **639** | **652** |
| | Opt.Rate | 100% | 100% | 100% | 100% |
| | Time(s) | 1.058 | 4.801 | 12.827 | 25.353 |

small-size instances denoted as 10*10, where 10 trains need to be dispatched at 10 stations. Subsequently, the learned model is allowed to dispatch the problem instances of a large scale. Its training process took about 6.25 hours, and Fig. 9 shows the averaged learning curve over five independent training processes. To test whether the RL model has learned dispatching skills, we compare the learned model with some handcrafted rules used in TTR, including First-Come-First-Service (FCFS) [19] and First-Schedule-First-Service (FSFS) [20]. Besides, we also compare our method with the exact algorithm implemented using Gurobi [38]. The performance of the learned model on 100 instances is shown in Table III.

In some cases, the FSFS strategy may generate infeasible solutions due to the track capacity constraints, so we approximate its performance by using the remaining feasible solutions. The proposed RL-based method with graph representation (RLG) outperforms all handcrafted methods, which means that the proposed method is effective in learning skills via trial and error. Meanwhile, the learned model, similar to the hand-made rules, fits well with the real-time requirements of the TTR task. Its scalability can also be proved by testing larger-scale problems, including 15 stations and 15 trains (15*15), 20 stations and 20 trains (20*20), 20 stations and 30 trains (20*30). With the help of the local search, the proposed method (RLG-LC) can get a significant improvement in solution quality with small computational costs. Besides, we find that Gurobi, as one of the state-of-the-art solvers, also presents an excellent performance with an acceptable computation time when facing train disturbances.

## C. TTR for Train Disruptions

With the help of the learning curriculum, our trained model can effectively address the TTR problem, even in cases of train disruptions. The associated learning curve is illustrated in Fig. 10. Table IV shows the experimental results of the learned model (named RGL) in terms of solution quality and computation time on 100 instances. To validate the effectiveness of the learning curriculum, we employ the vanilla PPO algorithm again to train our model in the complex scenarios, referred to as $RLG_L$, and also use the model learned in small-delay scenarios, denoted as $RLG_S$, directly in these complex scenarios. It is noteworthy that the performance of FSFS is not shown in this table, since it failed to find any feasible solution in almost all cases.

Regarding solution quality, we can observe that the learned model $RLG_S$ does not perform as well as it did under train disturbances, and $RLG_L$ performs comparably to the FCFS strategy, especially on larger-scale problems. By contrast, RGL achieves promising results, outperforming even the Gurobi solver in large-scale problems when using a similar computation time. With the help of the local search, its performance can be further improved, corresponding to RGL-LC. In terms of computational efficiency, these learned models can satisfy the real-time requirements of TTR in high-speed railways due to low computational costs. Even though the local search will increase the computation overhead to some extent, its response time is still acceptable for practical applications.

TABLE IV
THE TTR RESULTS UNDER TRAIN DISRUPTIONS

| Methods | Metrics | Test | | Generalization | |
|---|---|---|---|---|---|
| | | 10*10 | 15*15 | 20*20 | 20*30 |
| FCFS | Obj.(min) | 4792 | 10722 | 18407 | 24346 |
| | Gap | 16.24% | 35.10% | 46.26% | 61.83% |
| | Time (s) | **0.059** | **0.11** | **0.179** | **0.263** |
| $RLG_S$ | Obj. (min) | 6073 | 10416 | 14449 | 12869 |
| | Gap | 33.90% | 33.19% | 31.54% | 27.79% |
| | Time (s) | 0.123 | 0.282 | 0.469 | 0.863 |
| $RLG_L$ | Obj. (min) | 4711 | 10720 | 18780 | 24849 |
| | Gap | 14.80% | 35.08% | 47.33% | 62.60% |
| | Time (s) | 0.136 | 0.351 | 0.608 | 1.174 |
| RGL | Obj. (min) | 4260 | 7823 | 11758 | 12611 |
| | Gap | 5.77% | 11.04% | 15.87% | 26.31% |
| | Time (s) | 0.122 | 0.263 | 0.438 | 0.848 |
| RGL-LC | Obj. (min) | 4187 | **7533** | **11108** | **11576** |
| | Gap | 4.13% | 7.62% | 10.95% | 19.72% |
| | Time (s) | 2.53 | 4.188 | 5.245 | 5.897 |
| Gurobi | Obj. (min) | **4015** | 7629 | 15990 | 19874 |
| | Gap | 1.27% | 8.78% | 38.14% | 53.24% |
| | Time (s) | 1.27 | 5.937 | 8.264 | 10.214 |
| Baselines | Obj. (min) | 4014 | 6959 | 9892 | 9293 |
| | Opt.Rate | 100 % | 100 % | 100% | 98% |
| | Time (s) | 2.03 | 58.63 | 714 | 1116 |

## D. Ablation Studies

*1) **GNN-based state representation**:* To verify the effectiveness of the graph representations for TTR, we trained different models using the following methods: (1) using the state
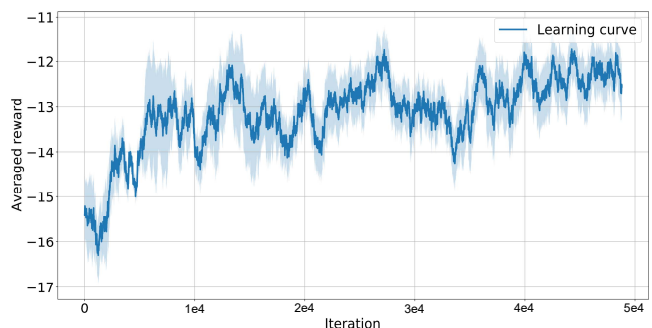


Fig. 10.  The learning curve of our model under train disruptions.

TABLE V
PERFORMANCE COMPARISON OF DECISION MODELS WITH DIFFERENT STATE REPRESENTATION IN TERMS OF TOTAL ARRIVAL DELAY $J$

| Problem size | $RLG_{raw}$ | $RLG_1$ | $RLG_2$ | $RLG_{AG}$ |
|---|---|---|---|---|
| 10*10 | 4991 | **4841** | 5264 | 5073 |
| 15*15 | 9849 | **8143** | 9675 | 9207 |
| 20*20 | 16852 | **12149** | 14831 | 13805 |
| 20*30 | 19540 | **13210** | 17463 | 14785 |

TABLE VI
PERFORMANCE COMPARISON OF DECISION MODELS WITH DIFFERENT LEARNING STRATEGIES IN TERMS OF TOTAL ARRIVAL DELAY $J$

| Problem size | $RLG_L$ | RGS | RGL | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | $w_4$=0.1 | $w_4$=0.5 | $w_4$=1 | $w_4$=1.5 | $w_4$=2 | $w_4$=5 | $w_4$=10 |
| 10*10 | 4711 | 4352 | 4379 | **4260** | 4282 | 4282 | 4301 | 4316 | 4427 |
| 15*15 | 10720 | 7974 | 8114 | 7823 | **7805** | 7849 | 7902 | 8126 | 8449 |
| 20*20 | 18780 | 12093 | 12641 | **11758** | 11896 | 11988 | 12128 | 12492 | 12944 |
| 20*30 | 24849 | 12997 | 15910 | **12611** | 12781 | 12895 | 12829 | 12734 | 12812 |

consisting of raw features in the graph, denoted as $RLG_{raw}$; (2) employing the extracted state from different numbers of the GNN layers, denoted as $RLG_n$, where $n$ represents the number of GNN layers. (3) using an off-the-shelf graph (i.e., the alternative graph (AG)) [39] to describe state, rather than the graph we designed, named $RLG_{AG}$. We evaluated the performance of these models in terms of total arrival delays $J$ of solutions, calculated by Eq. (1), and the experimental results on 100 instances are shown in Table V.

It can be found that the performance of the decision model is improved with the help of GNN layers. In comparison to the model using raw features (i.e., $RLG_{raw}$), the model with one GNN layer (i.e., $RLG_1$) consistently outperforms it on problems of all sizes. Although the performance of the decision model with two GNN layers (i.e., $RLG_2$) does not surpass that of $RLG_{raw}$ on the 10*10 problem, the influence of the GNN layer can still be observed on larger-scale problems. One possible reason is that when solving large-scale problems, the GNN can extract a more comprehensive perspective of the problem through the aggregation of node information. This global view of the problem may contribute to the decision model's ability to make better-informed decisions. Additionally, by comparing $RLG_1$ with $RLG_{AG}$, the superiority of our designed graph for TTR can be demonstrated. This is because we modify the direction of edges when the two nodes are from the same trains, such that the edge points from the latter event to the former one. As a result, the decision model can access future information, such as predicted train delays at following stations, when determining the priority between two trains at the current station.

*2) Learning Curriculum:* The proposed learning curriculum comprises two components, instance generation and loss definition. To assess their individual contributions, we initially trained our model using solely the sampling strategy, namely RGS. Subsequently, we examine the effect of the defined loss function by adjusting weight $w_4$ assigned to the loss item $L_4$. We present the experimental results on 100 instances in Table VI, where the performance is also defined as objective $J$.

As indicated in the table, the RGS model outperforms $RLG_L$, suggesting that our learning curriculum indeed enhances learning effectiveness. By selecting a value within the appropriate range of [0.5, 1.5] for hyperparameter $w_4$, the benefit of knowledge transfer can be found in the comparison between RGS and RGL.

*3) Local search:* In this part, we examine the influence of local search on our approach. Specifically, we vary the search
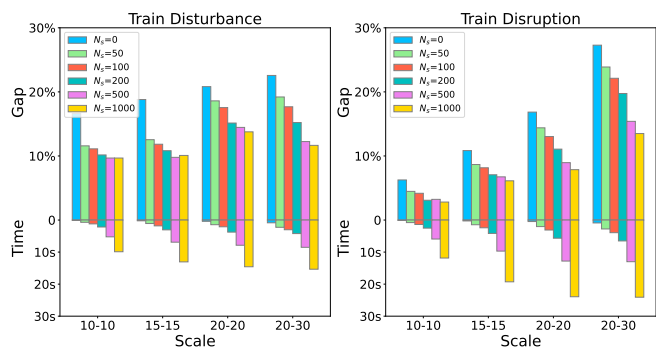


Fig. 11. A comparison of performance in terms of local search times.

times $N_s$ across different values, namely 0, 50, 100, 200, 500, and 1000. Fig. 11 shows the comparative results on problems with different scales, each evaluated on 100 instances.

The experimental results underscore the beneficial influence of employing local search to enhance the model's scalability. Specifically, when the learned model is applied in isolation ($N_s = 0$), its performance inevitably diminishes as the problem size increases, particularly in the case of train disruptions, but this issue can be alleviated to some degree by using local search. Especially, when the search times are set to 1000, the gap to the optimal solution can be reduced by up to 13%. Besides, the computation cost of local search is also acceptable for practical application. Even with 1000 search iterations, its computation time is less than 30 seconds. Furthermore, the advantages of local search become more apparent with the increase of the problem size or the severity of train delays. It enhances our confidence in applying our approach across diverse delay scenarios.

## VIII. CONCLUSION

In the daily operation of high-speed railways, providing a high-quality solution within a limited period is a grand challenge. In this study, we present a learning-based solution for TTR that directly satisfies the dispatcher's requirements for both response time and solution quality. In comparison to existing dispatching rules, the learned dispatching knowledge in our decision model can always obtain better performance across various problems with different scales and varying degrees of train delay, while maintaining reasonable computational costs. In particular, our method demonstrates superior performance over the Gurobi optimizer, a state-of-the-art solver, within a reasonable time frame when solving complex

scenarios involving train disruptions, highlighting the practical value of the proposed approach in these challenging scenarios. Moreover, our ablation studies demonstrate the benefits of the GNN-based state representation, learning curriculum and local search, further underscoring the potential of our approach in tackling the TTR problem.

However, it should be noted that there is still a significant amount of work that remains to be done. For instance, in cases of train disruptions, dispatchers may resort to more aggressive tactics to recover train delays, such as canceling trains and breaking connections. Although these decisions are infrequent, they are crucial, and this study has not taken them into account. A potential solution is to develop a multi-agent learning system in which each agent is responsible for a specific type of action, which will be our future work.

### REFERENCES

[1] P. Zhou, L. Chen, X. Dai, B. Li, and T. Chai, "Intelligent prediction of train delay changes and propagation using rvflns with improved transfer learning and ensemble learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 12, pp. 7432–7444, 2020.

[2] G. Cavone, T. van den Boom, L. Blenkers, M. Dotoli, C. Seatzu, and B. De Schutter, "An mpc-based rescheduling algorithm for disruptions and disturbances in large-scale railway networks," *IEEE Transactions on Automation Science and Engineering*, vol. 19, no. 1, pp. 99–112, 2020.

[3] V. Cacchiani, D. Huisman, M. Kidd, L. Kroon, P. Toth, L. Veelenturf, and J. Wagenaar, "An overview of recovery models and algorithms for real-time railway rescheduling," *Transportation Research Part B: Methodological*, vol. 63, pp. 15–37, 2014.

[4] R. Cheng, Y. Song, D. Chen, and L. Chen, "Intelligent Localization of a High-Speed Train Using LSSVM and the Online Sparse Optimization Approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 8, pp. 2071–2084, Aug 2017.

[5] P. Pellegrini, G. Marliere, R. Pesenti, and J. Rodriguez, "RECIFE-MILP: An Effective MILP-Based Heuristic for the Real-Time Railway Traffic Management Problem," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 5, pp. 2609–2619, Oct 2015.

[6] M. Dotoli, N. Epicoco, M. Falagario, B. Turchiano, G. Cavone, and A. Convertini, "A Decision Support System for real-time rescheduling of railways," in *2014 European Control Conference (ECC)*, Jun 2014, pp. 696–701.

[7] W. Fang, S. Yang, and X. Yao, "A survey on problem models and solution approaches to rescheduling in railway networks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 6, pp. 2997–3016, 2015.

[8] J. Rodriguez, "A constraint programming model for real-time train scheduling at junctions," *Transportation Research Part B: Methodological*, vol. 41, no. 2, pp. 231–245, 2007.

[9] K. Sato, K. Tamura, and N. Tomii, "A MIP-based timetable rescheduling formulation and algorithm minimizing further inconvenience to passengers," *Journal of Rail Transport Planning & Management*, vol. 3, no. 3, pp. 38–53, Aug 2013.

[10] B. Kersbergen, T. van den Boom, and B. De Schutter, "Distributed model predictive control for railway traffic management," *Transportation Research Part C: Emerging Technologies*, vol. 68, pp. 462–489, 2016.

[11] C. Zhang, Y. Gao, L. Yang, Z. Gao, and J. Qi, "Joint optimization of train scheduling and maintenance planning in a railway network: A heuristic algorithm using Lagrangian relaxation," *Transportation Research Part B: Methodological*, vol. 134, pp. 64–92, Apr 2020.

[12] Y.-H. Min, M.-J. Park, S.-P. Hong, and S.-H. Hong, "An appraisal of a column-generation-based algorithm for centralized train-conflict resolution on a metropolitan railway network," *Transportation Research Part B: Methodological*, vol. 45, no. 2, pp. 409–429, 2011.

[13] S. Zhan, S. Wong, P. Shang, Q. Peng, J. Xie, and S. Lo, "Integrated railway timetable rescheduling and dynamic passenger routing during a complete blockage," *Transportation Research Part B: Methodological*, vol. 143, pp. 86–123, Jan 2021.

[14] A. Bettinelli, A. Santini, and D. Vigo, "A real-time conflict solution algorithm for the train rescheduling problem," *Transportation Research Part B: Methodological*, vol. 106, pp. 237–265, Dec 2017.

[15] J. Törnquist Krasemann, "Design of an effective algorithm for fast response to the re-scheduling of railway traffic during disturbances," *Transportation Research Part C: Emerging Technologies*, vol. 20, no. 1, pp. 62–78, Feb 2012.

[16] S. Dündar and İ. Şahin, "Train re-scheduling with genetic algorithms and artificial neural networks for single-track railways," *Transportation Research Part C: Emerging Technologies*, vol. 27, pp. 1–15, 2013.

[17] S. Kanai, K. Shiina, S. Harada, and N. Tomii, "An optimal delay management algorithm from passengers' viewpoints considering the whole railway network," *Journal of Rail Transport Planning & Management*, vol. 1, no. 1, pp. 25–37, 2011.

[18] M. Wang, L. Wang, X. Xu, Y. Qin, and L. Qin, "Genetic algorithm-based particle swarm optimization approach to reschedule high-speed railway timetables: a case study in china," *Journal of Advanced Transportation*, vol. 2019, 2019.

[19] P. Wang, L. Ma, R. M. P. Goverde, and Q. Wang, "Rescheduling Trains Using Petri Nets and Heuristic Search," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 3, pp. 726–735, Mar 2016.

[20] W. Fang, S. Yang, and X. Yao, "A Survey on Problem Models and Solution Approaches to Rescheduling in Railway Networks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 6, pp. 2997–3016, Dec 2015.

[21] D. Šemrov, R. Marsetič, M. Žura, L. Todorovski, and A. Srdic, "Reinforcement learning approach for train rescheduling on a single-track railway," *Transportation Research Part B: Methodological*, vol. 86, pp. 250–267, Apr 2016.

[22] W. Li and S. Ni, "Train timetabling with the general learning environment and multi-agent deep reinforcement learning," *Transportation Research Part B: Methodological*, vol. 157, pp. 230–251, Mar 2022.

[23] H. Khadilkar, "A scalable reinforcement learning algorithm for scheduling railway lines," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 2, pp. 727–736, 2018.

[24] Y. Zhu, H. Wang, and R. M. Goverde, "Reinforcement Learning in Railway Timetable Rescheduling," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, 2020, pp. 1–6.

[25] Y. Wang, Y. Lv, J. Zhou, Z. Yuan, Q. Zhang, and M. Zhou, "A policy-based reinforcement learning approach for high-speed railway timetable rescheduling," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, 2021, pp. 2362–2367.

[26] L. Ning, Y. Li, M. Zhou, H. Song, and H. Dong, "A deep reinforcement learning approach to high-speed train timetable rescheduling under disturbances," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, 2019, pp. 3469–3474.

[27] P. Yue, Y. Jin, X. Dai, Z. Feng, and D. Cui, "Reinforcement learning for online dispatching policy in real-time train timetable rescheduling," *IEEE Transactions on Intelligent Transportation Systems*, 2023.

[28] W. Kool, H. Van Hoof, and M. Welling, "Attention, learn to solve routing problems!" *arXiv preprint arXiv:1803.08475*, 2018.

[29] I. Bello, H. Pham, Q. V. Le, M. Norouzi, and S. Bengio, "Neural Combinatorial Optimization with Reinforcement Learning," *5th International Conference on Learning Representations, ICLR 2017 - Workshop Track Proceedings*, pp. 1–15, nov 2016.

[30] A. Hottung, Y.-D. Kwon, and K. Tierney, "Efficient active search for combinatorial optimization problems," *arXiv preprint arXiv:2106.05126*, 2021.

[31] K. Zhang, F. He, Z. Zhang, X. Lin, and M. Li, "Multi-vehicle routing problems with soft time windows: A multi-agent reinforcement learning approach," *Transportation Research Part C: Emerging Technologies*, vol. 121, p. 102861, 2020.

[32] Y. Zhu, B. Mao, Y. Bai, and S. Chen, "A bi-level model for single-line rail timetable design with consideration of demand and capacity," *Transportation Research Part C: Emerging Technologies*, vol. 85, pp. 211–233, 2017.

[33] Y. Wang, Y. Lv, J. Zhou, Z. Yuan, Q. Zhang, and M. Zhou, "A policy-based reinforcement learning approach for high-speed railway timetable

rescheduling," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, 2021, pp. 2362–2367.

[34] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" *arXiv preprint arXiv:1810.00826*, 2018.

[35] C. Zhang, W. Song, Z. Cao, J. Zhang, P. S. Tan, and X. Chi, "Learning to dispatch for job shop scheduling via deep reinforcement learning," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1621–1632.

[36] Y. Wang, A. D'Ariano, J. Yin, L. Meng, T. Tang, and B. Ning, "Passenger demand oriented train scheduling and rolling stock circulation planning for an urban rail transit line," *Transportation Research Part B: Methodological*, vol. 118, pp. 193–227, 2018.

[37] P. R. d. O. da Costa, J. Rhuggenaath, Y. Zhang, and A. Akcay, "Learning 2-opt heuristics for the traveling salesman problem via deep reinforcement learning," in *Proceedings of The 12th Asian Conference on Machine Learning*, vol. 129, Nov 2020, pp. 465–480.

[38] "Gurobi optimizer reference manual," http://www.gurobi.com., 2022.

[39] F. Corman, A. D'Ariano, D. Pacciarelli, and M. Pranzo, "A tabu search algorithm for rerouting trains during rail operations," *Transportation Research Part B: Methodological*, vol. 44, no. 1, pp. 175–192, Jan 2010.