

# Hybrid intelligence for dynamic job-shop scheduling with deep reinforcement learning and attention mechanism

Yunhui Zeng<sup>a</sup>, Zijun Liao<sup>b</sup>, Yuanzhi Dai<sup>b</sup>, Rong Wang<sup>b</sup>, Xiu Li<sup>a</sup> and Bo Yuan<sup>a</sup>

<sup>a</sup>Shenzhen International Graduate School, Tsinghua University, Shenzhen 518055, China;

<sup>b</sup>School of Intelligent Systems Science and Engineering, Jinan University, Zhuhai 519070, China

## ARTICLE HISTORY

Compiled January 4, 2022

## Abstract

The dynamic job-shop scheduling problem (DJSP) is a class of scheduling tasks that specifically consider the inherent uncertainties such as changing order requirements and possible machine breakdown in realistic smart manufacturing settings. Since traditional methods cannot dynamically generate effective scheduling strategies in face of the disturbance of environments, we formulate the DJSP as a Markov decision process (MDP) to be tackled by reinforcement learning (RL). For this purpose, we propose a flexible hybrid framework that takes disjunctive graphs as states and a set of general dispatching rules as the action space with minimum prior domain knowledge. The attention mechanism is used as the graph representation learning (GRL) module for the feature extraction of states, and the double dueling deep Q-network with prioritized replay and noisy networks (D3QPN) is employed to map each state to the most appropriate dispatching rule. Furthermore, we present Gymjsp, a public benchmark based on the well-known OR-Library, to provide a standardized off-the-shelf facility for RL and DJSP research communities. Comprehensive experiments on various DJSP instances confirm that our proposed framework is superior to baseline algorithms with smaller makespan across all instances and provide empirical justification for the validity of the various components in the hybrid framework.

## KEYWORDS

DJSP, attention mechanism, reinforcement learning, hybrid intelligence, benchmark

## 1. Introduction

The job-shop scheduling problem (JSP) aims to determine the optimal sequential assignments of machines to jobs consisting of a series of operations subject to one or more criteria such as production cost and makespan. It is a typical combinatorial optimization problem closely related to industrial scheduling and production control (Satyro et al. 2021) and has been proven to be non-deterministic polynomial-time hard (NP-hard) (Garey, Johnson, and Sethi 1976) when the number of machines is greater than two. Furthermore, the emergence of smart manufacturing calls for effective solutions for the dynamic job-shop scheduling problem (DJSP) to confront the variabilities in production requirements, consumption of time, and available machines in the environment (Mohan, Lanka, and Rao 2019).

Among the existing solutions of DJSP, dispatching rules (Haupt 1989) such as Most Operation Remaining (MOR) and Shortest Processing Time (SPT) are designed based on practical experiences and can react to dynamic events promptly. However, there is a lack of assurance on their performance, and each of them is only suitable for specific scenarios. For example, MOR achieves lower makespan than SPT on ft06, an instance in the OR-Library (Beasley 1990), but SPT outperforms MOR on another instance la01. Consequently, the human operator needs to consistently choose one of the multiple commonly-used dispatching rules based on their own knowledge, which is impractical and tedious in most scenarios. Mathematical programming (Manne 1960) can find the optimal solution by solving an optimization problem with constraints. However, it suffers from the curse of dimensionality, and it is intractable to find the optimal scheduling in real-time, even for small-size instances. By contrast, meta-heuristic algorithms such as genetic algorithms (Gonçalves, Magalhães Mendes, and Resende 2005) and simulated annealing (Van Laarhoven, Aarts, and Lenstra 1992) aim to find a locally effective strategy by decomposing the DJSP into a series of static sub-problems and solving them separately (Luo 2020). The major issue is that these techniques feature poor generalization ability and need to be redesigned for even minor changes to the original problem.

With the development of artificial intelligence, data-driven approaches such as deep learning (Deng and Yu 2014) and machine learning (Das et al. 2015) have attracted great interest from researchers. Most of these methods rely on fully labeled datasets to train models and deploy them to unseen scheduling scenarios. With sufficient computing power and high-quality data sources, various encouraging results have been achieved following this supervised manner (Tian and Zhang 2021; Teymourifar et al. 2020). Unfortunately, for real-world production systems, it is often infeasible to obtain large-scale and high-quality training data.

Reinforcement learning (RL) aims at designing optimal policies to maximize the cumulative rewards over time through trial-and-error interactions with environments. It has achieved great success in many domains, such as GO (Silver et al. 2016; 2017), self-driving cars (Kendall et al. 2019), biology (Senior et al. 2018), surgical robotics (Richter, Orosco, and Yip 2019) and smart manufacturing (Xia et al. 2021). For the DJSP, Wang and Usher (2005) applied Q-learning to select self-designed dispatching rules on a single machine. The combination of deep learning and RL creates the field of deep reinforcement learning (DRL), which has been used to solve the JSP in recent years. Turgut and Bozdag (2020) applied the deep Q-network (DQN) to schedule jobs dynamically to minimize the delay time of jobs based on a discrete event simulation experiment. Wang et al. (2020) developed an improved Q-learning called dual Q-learning (D-Q) method for the assembly job shop scheduling problem (AJSSP). The top level Q-learning tries to learn the dispatching policy that can minimize machine idle rate, and the bottom level Q-learning tries to find the optimal scheduling policy to minimize the overall earliness of all jobs. Wang et al. (2021) employed proximal policy optimization (PPO) to alleviate the dimension disaster caused by the increase of problem size and obtain real-time production scheduling. Hameed and Schwung (2020) described a distributed RL method and applied deep deterministic policy gradient(DDPG) to each agent so that each agent optimizes its own cumulative reward.

A critical step in applying RL to the JSP is to formulate it as a Markov decision process (MDP), involving the key definitions of action, state and reward. However, previous studies often introduce customized formulations based on their own domain knowledge of JSP with handcrafted engineering and workload-specific implementation, making their experimental results incomparable. For the action space, some methods

select an action from a fixed operation set without any constraints (Turgut and Bozdag 2020; Luo et al. 2021; Qu, Wang and Shivani 2016), while others may select an action from the eligible operation set to ensure that the action can be directly executed (Zhang et al. 2020). There are also cases where self-designed (Wang et al. 2020; Wei and Zhao 2005) or customized dispatching rules (Lin et al. 2019; Yang and Yan 2009; Liu, Chang, and Tseng 2020) are designated as the actions. For the state space, some studies define the state as a matrix composed of customer order features and system features, such as the number of machines and the average completion time (Liu, Chang, and Tseng 2020), without specifying the relationship among different entities. Other studies use the disjunctive graph to represent the original state and the graph neural network (GNN) to do the feature extraction (Zhang et al. 2020). A major benefit is that they are capable of dealing with instances of varying sizes without extra training. However, the performance of the GNN may collapse and the computational cost may increase dramatically for complex problems (Wu et al. 2020), as the increasing number of neighborhood nodes may propagate noisy information (Zhou et al. 2020).

In summary, existing approaches deeply rely on the knowledge and experience of domain experts, which can be inevitably subjective and require significant implementation efforts. Meanwhile, there is a lack of common benchmarks and some studies even did not conduct experimental comparisons with other RL methods. Furthermore, the performance of RL-based scheduling systems is sensitive to the quality of formulation. Consequently, we argue that there is a critical need for a general, versatile and effective formulation scheme as well as a high-quality benchmark that is friendly to researchers and practitioners. Table 1 presents the details of existing RL-based JSP methods where FJSP, DFJSP, and FSP are the flexible job-shop scheduling problem, the dynamic, flexible job-shop scheduling problem, and the flow-shop scheduling problem, respectively.

Table 1. Existing RL-based methods for job shop scheduling.

Work	State representation	Action space	RL algorithm	Objective	Dynamic events	Problems	Comparison with other RL methods
Park et al. (2021)	GNN	Eligible operations	PPO	Makespan	None	JSP	Yes
Hamed and Schwung (2020)	GNN	Operations	PPO, DDPG	Makespan	None	JSP	No
Zhang et al. (2020)	GNN	Eligible operations	PPO	Makespan	None	JSP	No
Lin et al. (2019)	Matrix	Customized dispatching rules	multiclass DQN	Makespan	None	JSP	No
Wang et al. (2020)	Matrix	Self-designed rules	dual Q-learning	Total weighted earliness penalty, makespan	None	AJSP	Yes
Wang et al. (2021)	Matrix	Eligible operations	PPO	Makespan	Machine breakdown, Job rework	DJSP	No
Luo et al. (2021)	Matrix	Operations	DQN	Makespan	New job insertions	DJSP	Yes
Luo (2020)	Matrix	Self-designed rules	DQN	Tardiness	New job insertions	DFJSP	Yes
Turgut and Bozdag (2020)	Matrix	Operations	DQN	Job Delay time	New job insertions	DJSP	No
Yang and Yan (2009)	Matrix	Customized dispatching rules	B-Q learning	Tardiness	New job insertions	DJSP	No
Wei and Zhao(2005)	Matrix	Self-designed rules	Q-learning	Tardiness	Job delay	DJSP	No
Qu, Wang, and Shivani (2016)	Matrix	Operations	DQN	On-time delivery, the wait-in-process products	New job insertions	DJSP	No
Ariviv, Stern, and Edan (2016)	Matrix	Operations	dual Q-learning	Makespan	None	FSP	No
Reyna and Martínez-Jiménez, (2017)	Matrix	Operations	Q-learning	Makespan	None	FSP	No
Liu, Chang, and Tseng (2020)	Matrix	Customized dispatching rules	DDPG	Makespan	New job insertions	DJSP	Yes
Ours	Attention mechanism	General dispatching rules	D3QPN	Makespan	Machine breakdown, different order requirements	DJSP	Yes

In this paper, we propose a general framework to solve the DJSP based on the RL paradigm without handcrafted engineering. The output of RL for each state is a dispatching rule rather than detailed scheduling operations due to its clarity, ease of implementation, and good interpretability that are essential for the production field. In this setting, our framework can be regarded as a prototype towards hybrid intelligence, combining the formalized human knowledge in the form of structured decision rules and the data driven black-box intelligence exercised by the RL agent for selecting appropriate rules for the current state. The states of DJSP are represented by graphs and the attention mechanism is used as the graph representation learning (GRL) module, which is responsible for effective feature extraction. The extracted features are then used as the input to the double dueling deep Q-network with prioritized

replay and noisy networks (D3QPN), which is a value-based RL method with various useful extensions to DQN. To provide a principled DSJP environment for RL-based research, we develop a public benchmark named Gymjsp following the same convention as the OpenAI Gym environment that can significantly reduce the burden of other researchers. The overall procedure of the proposed framework is shown in Figure 1 and the major contributions of our work are summarised as follows:

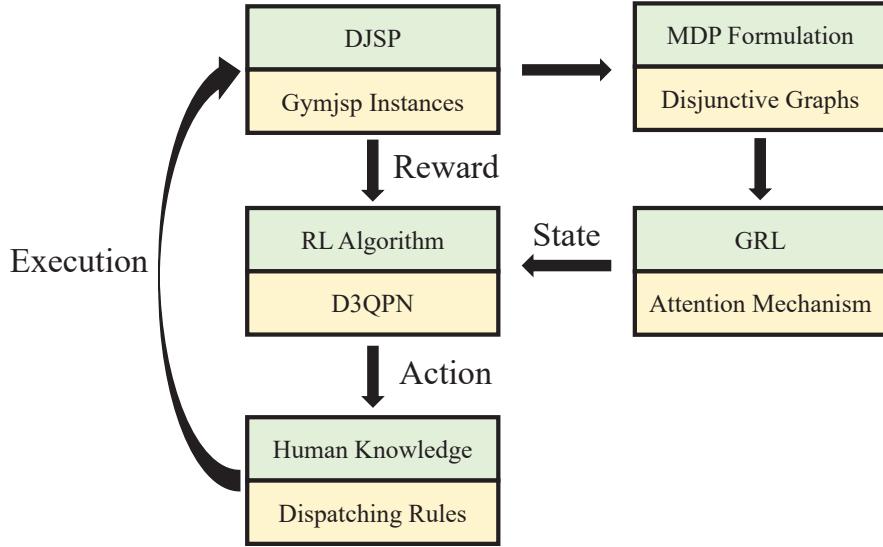


Figure 1. An overview of the proposed framework.

1. We effectively formulate the scheduling process of DJSP as an MDP, taking disjunctive graphs as states and general dispatching rules as the action. The disjunctive graph captures the local, global, and even dynamic information of JSP and the general dispatching rules can make our method applicable in a wide range of real smart manufacturing systems. Our formulation method provides a prototype towards hybrid intelligence, which can be an effective paradigm of human-machine cooperation.
2. We present a benchmark named Gymjsp following the convention of the OpenAI Gym environment to facilitate RL-based DJSP research. To the best of our knowledge, it is the first benchmark specifically targeted at RL-based DJSP.
3. We use the attention mechanism as the GRL module to extract features from disjunctive graphs. By visualizing the learning process of the attention mechanism, we confirm that the attention mechanism can effectively learn complex relationships among operations in the DJSP.
4. We propose an improved RL algorithm named D3QPN that has been proved to outperform other RL algorithms in the DJSP with several effective and complementary ingredients.

The rest of this paper is organized as follows: Section 2 provides a review of the relevant background, and the problem formulation is specified in Section 3. The details of our method are introduced in Section 4, and the experimental results and analyses are presented in Section 5. This paper is concluded in Section 6 with some discussions and directions for future work.

## 2. Background

### 2.1. Reinforcement learning

RL is one of the major areas of machine learning, along with supervised learning and unsupervised learning, which is inspired by biology and neuroscience. It is concerned about sequential decision making over time through trial-and-error interactions with the environment and the received reinforcement signals. In general, RL tasks are formulated as a MDP described by a tuple  $(\mathbf{S}, \mathbf{A}, \mathbf{P}, \mathbf{R})$ . At each decision point  $t$ , the agent will observe the state  $s_t \in \mathbf{S}$  and take a action  $a_t \in \mathbf{A}$  according to the policy  $\pi(\mathbf{S} \rightarrow \mathbf{A})$ , after which it enters a new state  $s_{t+1}$  with transition probability  $p(s_{t+1} | s_t, a_t) \in \mathbf{P}(\mathbf{S} \times \mathbf{A} \rightarrow \mathbf{S})$ . Meanwhile, an immediate reward signal  $r_t \in \mathbf{R}(\mathbf{S} \times \mathbf{A} \times \mathbf{S} \rightarrow \mathbb{R})$  is obtained as a result of the state transition  $(s_t, a_t, s_{t+1})$  (Sutton and Barto 2018). The goal of RL is to find an optimal policy  $\pi^*$  that maximizes the cumulative rewards:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E} \left[ \sum_{t=0}^{+\infty} r_{t+1} \mid \pi \right] \quad (1)$$

There are two main learning paradigms in RL: value-based methods and policy-based methods. A value-based method outputs the action-value of the pair  $(s, a)$ :  $Q^\pi(s, a) = \mathbb{E}^\pi \left[ \sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) \right]$ , where  $\gamma$  is the discount factor. By contrast, the policy-based method directly searches for an optimal policy using the gradient-based method and outputs the probability of each action. There is also a hybrid scheme called actor-critic where the actor interacts with the environment and updates the policy parameters, and the critic evaluates the policy's value to update the actor parameters, aiming to eliminate all the drawbacks from both value-based and policy-based methods.

### 2.2. Graph representation learning

A graph is a structure of objects (nodes or vertices), along with a set of relationships (edges or links) between each pair of objects. GRL aims to map a graph to a finite-dimensional vector space that captures the key features of the graph and assists with subsequent processing and decision making. One of the simplest methods to extract the key information from the graph is combining the features of nodes directly into a matrix, which can be viewed as a summarized graph feature. However, it cannot capture the relationship between each pair of nodes and requires manual feature filtering based on specific domain knowledge. Otherwise, the redundant information may be harmful to subsequent processing and decision making.

The GNN is a deep learning-based GRL model (Zhou et al. 2020) that has been widely applied in GRL recently (Cui et al. 2019; Chen et al. 2018), which can help nodes incorporate high-order neighborhood relationships (Shi et al. 2020). There is a hidden embedding  $\mathbf{h}_u^{(k)}$  corresponding to each node  $u$  in a GNN, which is updated to generate  $\mathbf{h}_u^{(k+1)}$  by aggregating the message  $\mathbf{m}_{\mathcal{N}(u)}$  from the graph neighborhood nodes  $\mathcal{N}(u)$  iteratively. The update function can be expressed as:

$$\begin{aligned} \mathbf{h}_u^{(k+1)} &= \mathbf{F}^{(k)} \left( \mathbf{h}_u^{(k)}, \mathbf{G}^{(k)} \left( \left\{ \mathbf{h}_v^{(k)}, \forall v \in \mathcal{N}(u) \right\} \right) \right) \\ &= \mathbf{F}^{(k)} \left( \mathbf{h}_u^{(k)}, \mathbf{m}_{\mathcal{N}(u)}^{(k)} \right) \end{aligned} \quad (2)$$

where both  $F$  and  $G$  are differentiable functions, while  $F$  is to update  $\mathbf{h}_u^{(k)}$  and  $G$  is to aggregate message  $\mathbf{m}_{\mathcal{N}(u)}$ . After  $K$  iterations, the embeddings of each node can be defined using the outputs of the final layer  $\mathbf{z}_u = \mathbf{h}_u^{(K)}$ ,  $\forall u \in \mathcal{V}$  (Hamilton 2020). Note that, the computation on  $\mathcal{N}(u)$  can lead to dramatically increased time consumption given a complex graph, resulting in deteriorated performance.

The attention mechanism can also be extended to GRL. Its main idea is to aggregate a variable-sized input by a function while focusing on the most relevant parts of the input to make decisions. Each node in the set of neighborhood nodes  $N(t)$  of the target node  $t$  is assigned an attention weight  $\alpha_i$  using the attention-based GRL method, where  $\alpha$  is from 0 to 1 with  $\sum_{i \in N(t)} \alpha_i = 1$ . The attention mechanism can identify the information in the input graph pertinent to subsequent tasks by calculating the importance of different nodes. Furthermore, many operations in attention can be performed in parallel, greatly improving computational efficiency.

### 2.3. *Simulation environment for DJSP*

The research on JSP has significantly benefited from several simulation environments containing standard instance sets for simulating real-world scheduling environments and evaluating different algorithms. For example, Demirkol's instances (Demirkol, Mehta, and Uzsoy 1998) are for the problems of maximum lateness and minimum makespan in job shops and flow shops. Taillard (1993) created instances for the job shop, open shop, and permutation flow shop scheduling problems with the objective of minimizing the makespan. In addition, the well-known instance set, OR-Library, has been widely adopted as the test environment of JSP in recent years. These instances contain the key parameters and data of JSP, such as the processing time of each operation, the number of jobs, the number of machines, etc. However, these simulation environments were originally designed for traditional methods such as dispatching rules, which cannot be directly used for RL-based approaches. Consequently, researchers need to adapt these environments accordingly to suit RL methods, which is challenging and time-consuming. Furthermore, it is common that researchers tend to repackage the environments based on their own domain knowledge and specific purposes. For example, the information used to define the state is usually different: some researchers take into account the consumption of time and status of operations (Zhao and Zhang 2021), while others may consider the status of jobs and machines when designing their states (Wang et al. 2021). For the DJSP, researchers may set different effects (e.g., the length of time delays) for the same dynamic event (e.g., machine breakdown), which makes the experimental results incomparable. In summary, there is a lack of a public, easy-to-use, and standardized benchmark that can help promote the development of RL in the DJSP.

## 3. Problem formulation

### 3.1. *DJSP formulation*

The DJSP is a dynamic allocation problem in a resource-limited environment facing variability in production requirements, consumption of time, available machines and unexpected faults. In each DJSP, there are  $m$  machines  $M = \{M_1, M_2, \dots, M_m\}$  and  $n$  jobs  $J = \{J_1, J_2, \dots, J_n\}$ . Each job has  $m$  operations to be processed where the operation sequence of  $J_i$  is defined as  $O_i = \{O_{i1}, O_{i2}, \dots, O_{il}, \dots, O_{im}\}$ . The machine matrix  $MO =$

$\{M_{il} \mid M_{il} = 0, M_1, M_2, \dots, M_m\}$  ( $i = 1, 2, \dots, n, l = 1, 2, \dots, m$ ) indicates that  $O_{il}$  should be processed by  $M_{il}$ . The value of  $T_{il}$  in the processing time matrix  $TO = \{T_{il} \mid T_{il} \geq 0\}$  ( $i = 1, 2, \dots, n, l = 1, 2, \dots, m$ ) specifies the required processing time of  $O_{il}$ . If  $M_{il} = 0$ , then  $T_{il} = 0$ , which indicates that job  $J_i$  does not require the operation  $O_{il}$  (Wang et al. 2021).

In our work, the goal of DJSP is to minimize the makespan with the following objective function:

$$\min \left\{ \max_{1 \leq i \leq n} C_{il} \right\} \quad (3)$$

where  $C_{il}$  is the completion time of operation  $O_{il}$  on machine  $m_{il}$ . As the machine matrix  $MO$  and the processing time matrix  $TO$  are given in advance, the completion time  $C_{il_{il}}$  can be given as:

$$\begin{aligned} C_{il} &= \max(C_{i(l-1)}, C_{ab}) + T_{wait} + T_{il} \\ M_{il} &= M_{ab}; i, a = 1, \dots, n; b = 1, \dots, m; l = 2, \dots, m \end{aligned} \quad (4)$$

where  $C_{i(l-1)}$  is the completion time of the previous operation  $O_{i(l-1)}$  in the same job  $J_i$ ;  $C_{ab}$  is the completion time of the previous operation  $O_{ab}$  on the same machine  $M_{il}$ ;  $T_{wait}$  indicates the time that it takes to keep machine  $M_{il}$  waiting without loading operation  $O_{il}$  when the machine  $M_{il}$  is idle, which is ineffective in most cases.

The entire process of a simple DJSP example is shown in Figure 2, where machines in red have the possibility of breakdown. From the circular starting point to the square ending point, each job is processed by a certain sequence of machines, resulting in a unique route of completion. The simultaneous execution of multiple jobs may inevitably lead to competition for machine resources, blocking the execution of jobs and increasing the total time cost.

The assumptions in the DJSP are listed as follow:

1. Each machine can conduct only one operation at a time.
2. Each operation of a job can be executed by only one machine at a time.
3. Each operation of a job cannot be executed until its primary operations are completed.
4. Each operation that has been started cannot be suspended or terminated.
5. All the operations of the same job need to follow a unique sequence.
6. All the transportation times and setup times are negligible.
7. The machine matrix  $MO$  and processing time matrix  $TO$  are known in advance.
8. There is some possible disturbance in the production environment, such as machine breakdown.

### 3.2. MDP formulation for DJSP

The DJSP aims to determine the order of operations in a dynamic environment, which is essentially a sequential decision problem and can be formulated as an MDP defined by states, actions, and reward functions.

States describe a common feature set that can represent the characteristics related to the environment and the objective of scheduling. We use the disjunctive graph

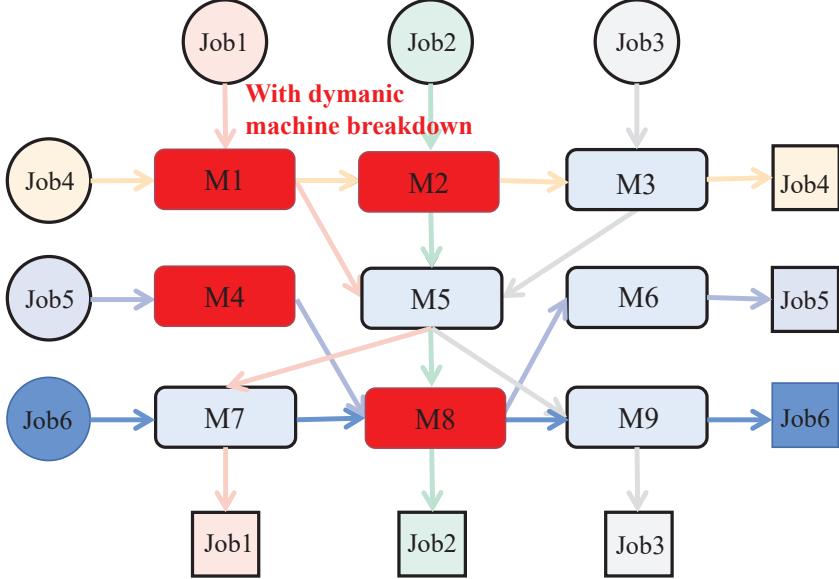


Figure 2. A simple example of DJSP with nine machines and six jobs where the arrows of the same color indicate the unique flow of a specific job.

whose details can be found in Section 3.3 as the abstract of the local, global, and even dynamic information of the scheduling environment. Furthermore, a disjunctive graph with a numerical representation of the state attributes is capable of specifying all states of different scheduling problems. In practice, we can rely on domain knowledge to extract features from disjunctive graphs manually or use GRL methods such as the GNN for automatic feature extraction.

Traditionally, the action at each step is a detailed operation to be directly executed in the environment, which can be unsafe for RL-based methods due to the black-box nature of neural networks. Instead, we adopt dispatching rules as the action space of our RL method for the sake of simplicity, ease of implementation, and good interpretability, which are desirable for management purposes. As an analogy to the factory staff who routinely choose one of the dispatching rules based on their domain knowledge, our RL-based method strategically selects the most appropriate dispatching rule according to the current state. In this way, it is possible to overcome the limitation of any single dispatching rule and can be seen as a hybrid intelligence method. Figure 3 shows the execution flows of eight dispatching rules adopted in our work. Note that, in addition to rule-based systems, we can use any scheduling methods as possible actions such as meta-heuristic and RL-based algorithms.

The definition of reward should be closely associated with the scheduling objective. Although the goal of the DJSP is to minimize the makespan, it can only be obtained when the entire scheduling process is over, resulting in the challenging sparse reward issue (Rauber et al. 2021). Since the makespan is closely related to machine utilization (in general, the higher the machine utilization, the smaller the makespan), our reward function is defined as:

$$reward = -\frac{\sum_1^k \frac{\text{number of idle machines}}{\text{number of total machines}}}{k} \quad (5)$$

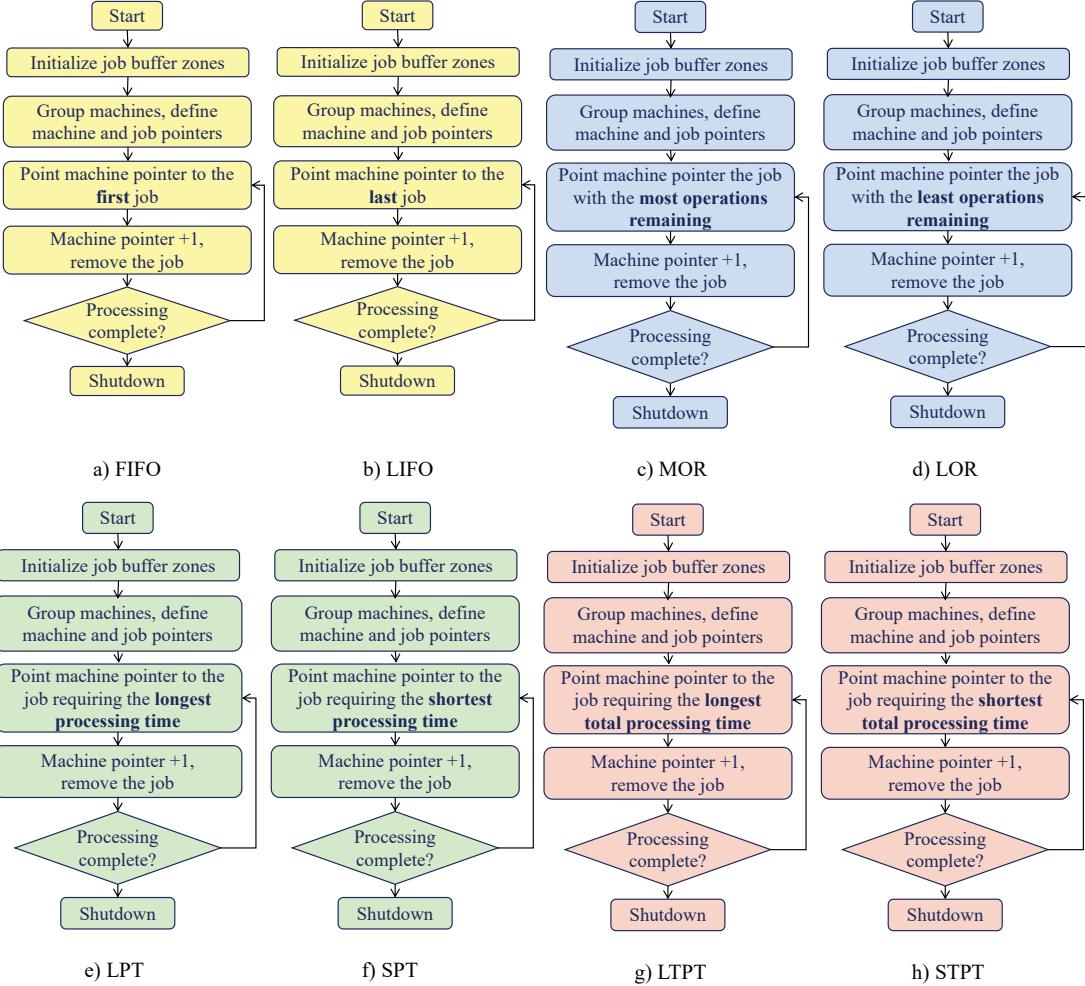


Figure 3. The flowcharts of eight rule-based scheduling algorithms: (a) First In First Out (FIFO); (b) Last In Last Out (LIFO); (c) Most Operations Remaining (MOR); (d) Least Operations Remaining (LOR); (e) Longest Processing Time (LPT); (f) Shortest Processing Time (SPT); (g) Longest Total Processing Time (LTPT); (h) Shortest Total Processing Time (STPT).

where  $k$  is a hyperparameter controlling the number of steps for which the output dispatch rule is used. Although it is possible to design the reward function in more sophisticated ways, our purpose is to show that, with such a simple reward function, the proposed hybrid RL framework can already achieve superior performance, with plenty of room for further improvement.

### 3.3. State representation

The performance of RL on the JSP relies heavily on the quality of hand-crafted feature representation. The disjunctive graph  $G = (V, C \cup D)$  is a popular representation for JSP instances (Balas 1969).  $V$  is the set of all possible operations represented as vertices in the graph, including two dummy vertices: the start and terminal ones, both

of which have zero consumption of time.  $C$  is a set of directed edges (conjunctions) representing the precedence constraint between two consecutive operations in the same job.  $D$  is the undirected edge set (disjunctions) representing the machine-sharing constraint between two vertices. When a pair of operations can be processed on the same machine, their corresponding vertices are connected with an undirected edge. Therefore, solving a JSP instance is equivalent to turning every undirected disjunctive edge into a directed edge so that the result is a directed acyclic graph. Figure 4 gives an example of a simple JSP instance and its solution.

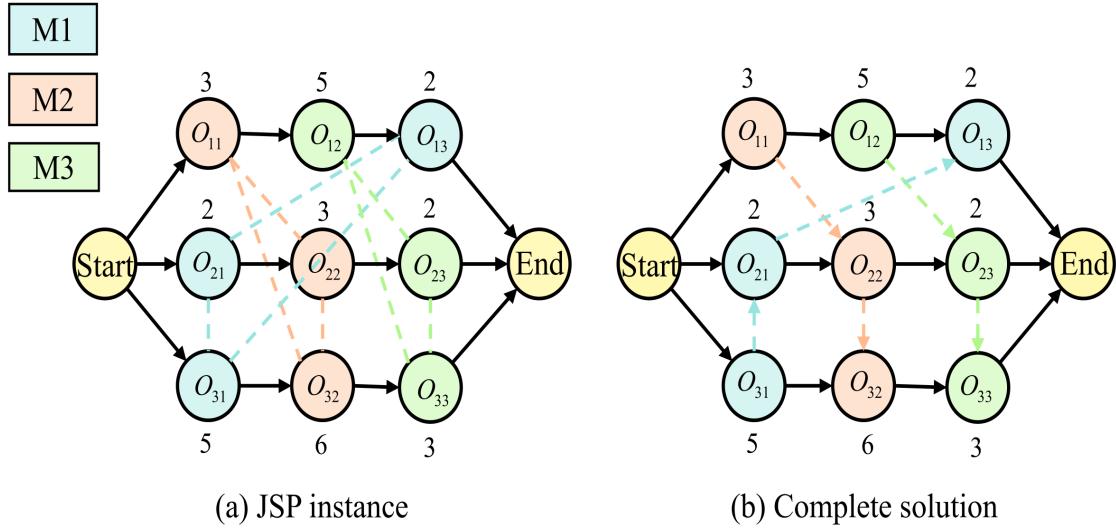


Figure 4. An example of a  $3 \times 3$  JSP instance. (a) Initial state. Each circle is an operation and  $O_{ij}$  is the  $j^{th}$  operation of the  $i^{th}$  job. The number on each circle represents the consumption of time of the operation. The solid and dotted lines are conjunctive and disjunctive edges, respectively. Vertices of the same color indicate that the operations are processed on the same machine; (b) Its complete solution.

Although the disjunctive graph contains the static information in the JSP, such as precedence constraints and machine-sharing constraints, it fails to represent the dynamic state information in the DJSP. To address this issue, we add several features into the vertices in  $G$  as follows:

1. Job ID: The sequence numbers of jobs that contain this operation.
2. Step ID (Operation ID): The sequence number of this operation.
3. Node type: -1 for unfinished, 0 for in the process, and 1 for completed.
4. Completion ratio: The completion rate of the whole job when this operation is completed.
5. Consumption of time: The required time to complete the operation.
6. Number of remaining operations: The number of subsequent operations.
7. Waiting time: The time elapsed from the beginning until this operation can be processed.
8. Remaining time: The remaining consumption of time (0 for not in processing).
9. Doable: *True* if this operation is doable.
10. Machine ID: The sequence number of the machine that can process this operation (0 if this operation cannot be processed).

## 4. Methodology

### 4.1. Proposed architecture for DJSP

This paper proposes a hybrid framework to solve the DJSP in smart manufacturing using the attention mechanism and RL, as shown in Algorithm 1. The attention mechanism is a GRL module to transfer the original state  $s$ , which is a disjunctive graph, to the extracted state  $s_a$ , which is a vector, and D3QPN is used to learn the action value  $Q(s_a, \cdot, \theta)$ . In Algorithm 1, given a DJSP, the system initially formulate the DJSP as an MDP ( $\mathbf{S}, \mathbf{A}, \mathbf{P}, \mathbf{R}$ ), where  $\mathbf{S}$  is the disjunctive graph (Line 2). Then, the attention mechanism transfers this disjunctive graph to a vector that can be directly input to D3QPN (Line 4). With the sampled noisy network  $\xi$  (Line 5) available, we retrieve the dispatching rule corresponding to the maximal value (Line 6). This dispatching rule is executed for  $k$  steps to make the control process more stable (Line 7 to Line 9) and the transition  $(\mathbf{s}_a, \mathbf{a}, \mathbf{r}, \mathbf{s}'_a)$  is stored in the replay buffer with maximal priority to make sure that each experience is seen at least once. The RL algorithm and GRL parameters are updated using the mini-batch sampled from the replay buffer. The above procedure is repeated until the maximal number of epochs  $max\_epoch$  is reached.

---

#### Algorithm 1 Proposed framework using D3QPN and attention mechanism

---

**Input:** Environment  $Env$ , set of random variables of the network  $\epsilon$

- 1: Initialize the noisy behavior network  $Q$  with random weights  $\theta$ ; GRL module; replay buffer; exponent  $\alpha$  of prioritized replay; minibatch  $\tau$ .
- 2: Formulate the DJSP as an MDP( $\mathbf{S}, \mathbf{A}, \mathbf{P}, \mathbf{R}$ ), where  $\mathbf{S}$  is a disjunctive graph.
- 3: **for** epoch number  $epoch = 1, 2, \dots max\_epoch$  **do**
- 4:   Extract the state  $s_a$  from  $s$  using the GRL module (Algorithm 2).
- 5:   Sample a noisy network  $\xi \sim \epsilon$ .
- 6:   Select the dispatching rule  $a \leftarrow \text{argmax}_{a_t \in A} Q(s_a, \cdot, \xi, \theta)$ .
- 7:   **for** schedule cycle  $i = 1, 2, \dots k$  **do**
- 8:     Execute dispatching rule  $a$  and observe new disjunctive graph  $s'$ .
- 9:     Extract the state  $s'_a$  from  $s'$  using Algorithm 2.
- 10:   **end for**
- 11:   Store  $(\mathbf{s}_a, \mathbf{a}, \mathbf{r}, \mathbf{s}'_a)$  in replay buffer with maximal priority  $p_t = \max_{i < t} p_i$ .
- 12:   Sample the minibatch of transitions with probability  $Per(j) = p_j^\alpha / \sum_\tau p_\tau^\alpha$ .
- 13:   Update the learning algorithm and GRL module based on Algorithm 3.
- 14: **end for**

**Output:** The learned RL module and the GRL module

---

### 4.2. Graph representation learning using attention mechanism

Our GRL module uses stacked self-attention, residual connection (He et al. 2016), layer normalization (LN) (Ba, Kiros, and Hinton 2016) and the fully-connected feed-forward network (FFN) to map the sequence that contains all operations  $X^0 = (x_{11}^0, x_{12}^0, \dots, x_{nm-1}^0, x_{nm}^0)$  in the input disjunctive graph to a sequence of extracted features  $X^L = (x_{11}^L, x_{12}^L, \dots, x_{nm-1}^L, x_{nm}^L)$ , which has the same shape as the input sequence, where  $L$  is the number of GRL module layers. The overall structure of the proposed GRL module is shown in Figure 5.

Self-attention is a function that maps a query and a set of key-value pairs to an output. It learns three linear projections to map the state to  $Q, K, V \in \mathbb{R}^{d_{length} \times d_{feature}}$ ,

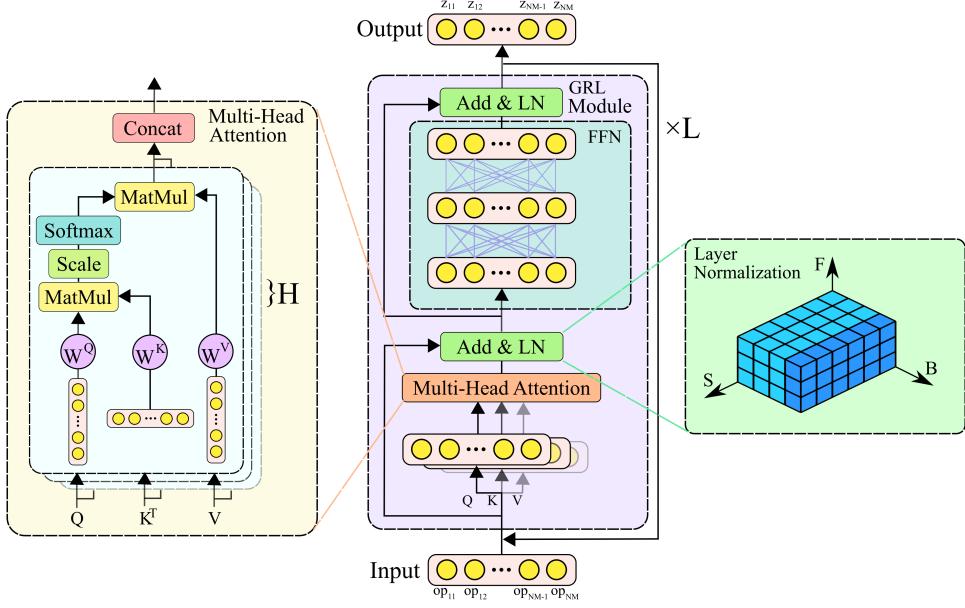


Figure 5. Overall structure of the proposed graph representation learning module.

where  $d_{length}$  is the number of total operations, and  $d_{feature}$  is the dimension of the features. The query  $Q$  indicates the type of operations that the current operation concerns; the key  $K$  indicates the type of the current operation; the value  $V$  contains the information of this operation. In the first step, we compute the matching scores among all operations by the matrix multiplication of queries and keys. The scores are then divided by  $\sqrt{d_{feature}}$  to prevent the attention from being over-focused on nodes with large scores (Vaswani et al. 2017). Next, we apply a softmax function to obtain the attention weight on the values. Finally, we conduct the matrix multiplication of the attention weight and value to obtain the  $output \in \mathbb{R}^{d_{length} \times d_{feature}}$ . The aforementioned processes can be executed in parallel, which can greatly improve computational efficiency. The procedure of the attention function is as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK'}{\sqrt{d_{feature}}}\right)V \quad (6)$$

Multi-head attention acquires more information by integrating the outputs of multiple attention functions with different query matrices and key matrices.

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_H) W^{multi} \\ \text{where } \text{head}_i &= \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right) \end{aligned} \quad (7)$$

where  $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{feature} \times \frac{d_{feature}}{H}}$  and  $W^{multi} \in \mathbb{R}^{d_{feature} \times d_{feature}}$  are the parameter matrices, and  $H$  is the number of attention heads.

Following the multi-head attention layer, the FFN is composed of two layers of linear transformation, which maps the input from  $d_{input}$  dimensions to  $d_{output}$  dimensions

using the parameter  $W_1, W_2, b_1$  and  $b_2$ . We also apply residual connections around each of the multi-head attention and the FFN to ease the issue of gradient disappearance. After the residual connection, we employ LN, a regularization method, to avoid internal covariate shifts when there are few examples in a single batch. Algorithm 2 gives the complete process of the proposed GRL.

---

**Algorithm 2** Graph representation learning using attention mechanism

---

**Input:**  $X^0 = \{x_{11}^0, \dots, x_{nm}^0\}$ : Disjunctive features

```

1: Initialize the following parameters:
2: (1) $L$ : number of GRL module layers;
3: (2) $H$ : number of head of Multi-Head Attention;
4: (3) $W^{Q^l}, W^{K^l}, W^{V^l}, W^{multi^l}, W_1^l, W_2^l$ : weight matrix, where  $l = 1, \dots, L$ ;
5: (4)  $\gamma_1^l, \gamma_2^l, \beta_1^l, \beta_2^l$ : parameters for LN function;
6: (5)  $b_1^l, b_2^l$ : bias for FFN.
7: for  $l \leftarrow 1$  to  $L$  do
8:   for  $h \leftarrow 1$  to  $H$  do
9:      $Q^l \leftarrow X^{l-1} \cdot W^{Q^l}, K^l \leftarrow X^{l-1} \cdot W^{K^l}, V^l \leftarrow X^{l-1} \cdot W^{V^l}$ .
10:     $head_h^l \leftarrow \text{softmax}\left(\frac{Q^l \cdot K^l}{\sqrt{d_{feature}}}\right) V^l$ .
11:   end for
12:    $A^l \leftarrow \text{Concat}(head_1^l, \dots, head_H^l) \cdot W^{multi^l}$ .
13:    $X^l \leftarrow X^{l-1} + A^l$ .
14:    $X^l \leftarrow \text{LN}_1^l(X^l)$ .
15:    $F^l \leftarrow \max(0, X^l \cdot W_1^l + b_1^l) \cdot W_2^l + b_2^l$ .
16:    $X^l \leftarrow X^l + F^l$ .
17:    $X^l \leftarrow \text{LN}_2^l(X^l)$ .
18: end for
```

**Output:**  $X^L = \{x_{11}^L, \dots, x_{nm}^L\}$ : Extracted feature

---

### 4.3. Double dueling DQN with prioritized replay and noisy networks

The DQN is an important milestone in RL that can achieve competitive performance compared to humans on various tasks. At each step of DQN, the agent selects an action according to the  $\epsilon$ -greedy strategy with respect to the action value and adds a transition into the replay buffer. Then, the agent uses a sampled batch from the replay buffer to optimize the behavior neural network to minimize the loss:

$$\left( R_{t+1} + \gamma \max_{a'} Q_{\bar{\theta}}(s_{t+1}, a_{t+1}) - Q_{\theta}(s_t, a_t) \right)^2 \quad (8)$$

where  $\gamma$  is a discount factor. The behavior network is updated every step, while the target network is updated every specific number of steps to ensure the stability of the action value. In this paper, we propose a double dueling DQN with prioritized replay and noisy networks (D3QPN) algorithm (Algorithm 3), which combines four effective and complementary ingredients with DQN to address the limitations of DQN:

1. Double Q-learning (Van Hasselt, Guez, and Silver 2016), which is used to reduce harmful overestimations of action value by decoupling the action values used to

---

**Algorithm 3** The training procedure of D3QPN

---

**Input:** step-size  $\eta$ , target network replacement frequency  $N^-$ , exponent  $\beta$  of prioritized replay, replay size  $N$

- 1: Initialize noisy behavior network  $Q$  with random weights  $\theta$ ;
- 2: Initialize target network  $\hat{Q}$  with weights  $\theta^- = \theta$ ;
- 3: Initialize the GRL module  $G(\phi)$ .
- 4: **for**  $j=1$  to  $\tau$  **do**
- 5:   Sample noisy variables  $\xi' \sim \varepsilon$  for target network .
- 6:   Sample noisy variables  $\xi'' \sim \varepsilon$  for behavior network.
- 7:   Compute importance-sampling weight  $w_j = (N \cdot Per(j))^{-\beta} / \max_i w_i$ .
- 8:   Set  $y_j = \begin{cases} r_j & \text{for terminal} \\ r_j + \gamma \max_{a'} Q(s_j, \arg \max_{b \in \mathcal{A}} Q(y_j, b, \xi''; \theta), \xi'; \theta') & \text{for non-terminal} \end{cases}$
- 9:   Compute TD-error  $\delta_j = (y_j - Q(s_j, a_j; \xi; \theta))^2$ .
- 10:   Update transition priority  $p_j \leftarrow |\delta_j|$ .
- 11:   Accumulate weight-change  $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_\theta Q(s_j, a_j)$ .
- 12: **end for**
- 13: Update weights  $\phi \leftarrow \phi + \eta \cdot \Delta$ ,  $\theta \leftarrow \theta + \eta \cdot \Delta$ , reset  $\Delta = 0$ .
- 14: Every  $N^-$  times, update target network:  $\theta^- \leftarrow \theta$ .

**Output:**  $Q(\cdot, \varepsilon; \theta)$  action-value function, GRL module  $G(\phi)$

---

select and to evaluate an action with the loss:

$$\left( R_{t+1} + \gamma_{t+1} Q_{\bar{\theta}} \left( s_{t+1}, \arg \max_{a'} q_{\theta} (s_{t+1}, a') \right) - Q_{\theta} (s_t, a_t) \right)^2 \quad (9)$$

2. Prioritized replay (Schaul et al. 2015), which is used to sample transitions with probability  $p_t$  relative to TD error, making the transition that contains more information has a higher probability of being sampled. New transitions are given maximum priority when inserted into the replay buffer to ensure that all transitions are used at least once.
3. Dueling networks (Wang et al. 2016), which have two streams to separately estimate the state-value function and the action advantage function for each action and can aggregate the state-action value function without imposing any change to the preceding RL algorithm:

$$q(s_t, \mathbf{a}_t) = v(s_t) + A(s_t, \mathbf{a}_t) \quad (10)$$

4. Noisy networks (Fortunato 2017), an effective exploration and exploitation strategy using the factorised Gaussian noise. The parameters of the noisy layer are:

$$\begin{aligned} w &= \mu^w + \sigma^w \odot \varepsilon^w \\ b &= \mu^b + \sigma^b \odot \varepsilon^b \end{aligned} \quad (11)$$

where the  $\mu^w, \mu^b, \sigma^w$  and  $\sigma^b$  are the parameters of neural network;  $\varepsilon^b$  and  $\varepsilon^w$  are the random variables of factorised Gaussian noise;  $\odot$  denotes the element-wise product. The output of the noisy layer is  $y = wx + b$ , which is the same as the standard fully connected layer.

The architecture of D3QPN is shown in Figure 6, which integrates all the components mentioned above into a single agent.

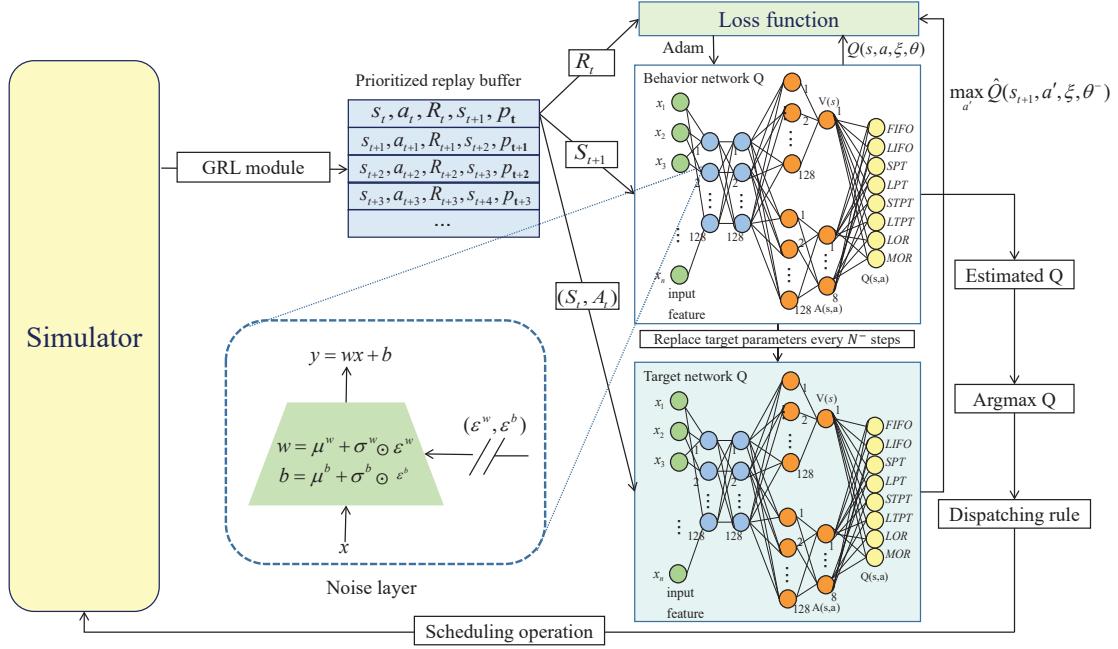


Figure 6. The double dueling DQN with prioritized replay and noisy networks.

## 5. Experiments

### 5.1. Gymjsp

As mentioned in Section 2.3, there is a critical lack of standardized practice in current RL-based JSP research in terms of experimental simulators and evaluation procedures, making it difficult for algorithm comparison. Furthermore, existing simulators are not friendly to researchers to implement their RL algorithms on JSP, which may seriously hinder the development of RL-based JSP solutions. To bridge this gap, we propose a new public benchmark named **Gymjssp**<sup>1</sup> based on the well-known OR-Library, a collection of datasets for a variety of Operations Research (OR) problems including JSP, portfolio optimization and so on. To provide an easy-to-use suite for RL algorithms, Gymjsp features an interface similar to OpenAI Gym (Brockman et al. 2016), which is a widely used benchmark in the RL community. According to Figure 7, the scheduling environment is initialized based on the instances of the OR-Library and, if required, with disturbance on the time consumption (a longer time consumption than normal one of an operation on a machine can be viewed as the event of machine breakdown for the next operation). It is also possible to control whether to disorder the processing sequence of operations by the parameter **shuffle**, representing different order requirements. All the dynamic effects can be reproduced using the same random seeds by the function **seed**. During initialization, with the parameter **random rate**, the time consumptions of all operations are determined as follows:

<sup>1</sup><https://github.com/Yunhui1998/Gymjsp>

$$T'_{i,l} = T_{i,l} + \text{noisy}_{(i,l)} \quad (12)$$

$$\text{noisy}_{(i,l)} = \begin{cases} \min(1, \max(-1, \mathcal{N}(0, 0.1))) \cdot T_{i,l}, & r < \text{random\_rate} \\ 0, & \text{other} \end{cases} \quad (13)$$

where  $T'_{i,j}$  is the new time consumption for the  $j^{th}$  operation of the job  $J_i$ ;  $r \in [0, 1]$  is a random float;  $\text{random\_rate} \in [0, 1]$  is a hyperparameter.

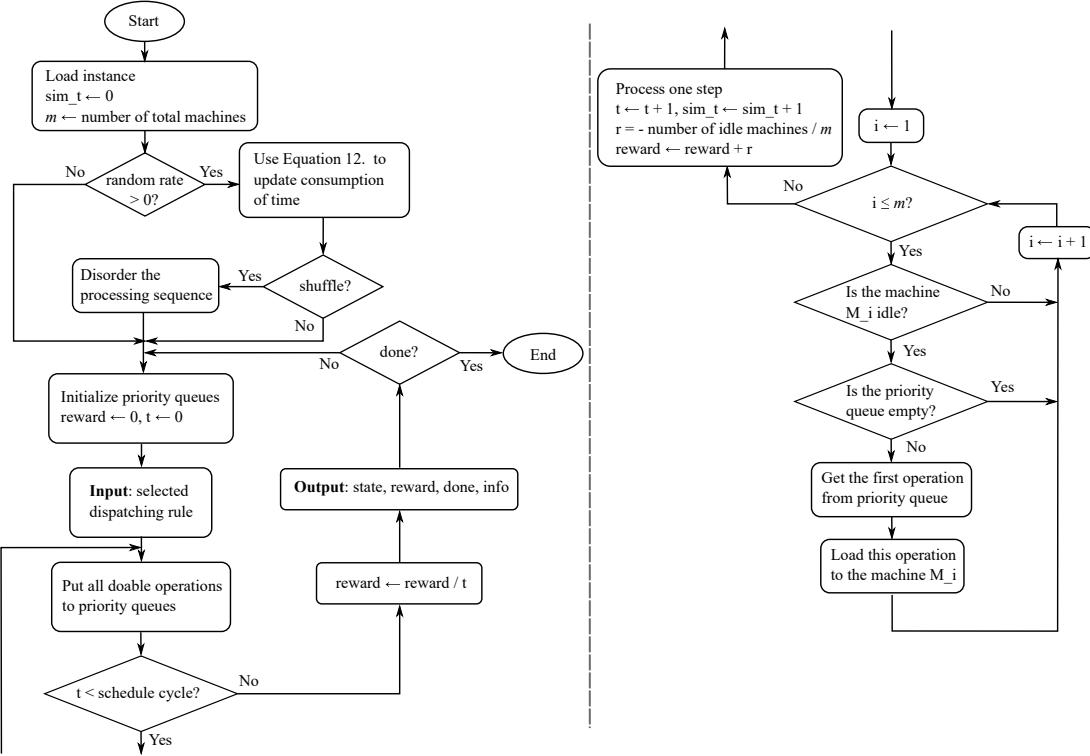


Figure 7. The scheduling process based on Gymjsp. **sim\_t**: the simulated global time; **t**: the accumulated execution time of the current rule; **done**: the indicator of whether all jobs have finished; **info**: extra information such as makespan (zero if there are still unfinished operations).

After environment initialization, in each decision step, a dispatching rule from {FIFO, LIFO, LPT, SPT, STPT, LPTP, MOR, LOR} is selected as the action. Each dispatching rule assigns a specific priority queue for each machine and, for machine  $M_i$ , its priority queue  $q_i$  contains, in ascending order, all doable operations to be processed by  $M_i$ . Next, the environment schedules operations following this rule for **schedule cycle** times and returns the average ratio of idle machines as the reward of this action. Finally, when all the jobs are finished, the **done** flag is set to **True**, and the extra information such as makespan is stored in **info** for debugging or testing.

## 5.2. Performance evaluation

The key objective of the experimental studies is to investigate the advantage of our proposed framework of hybrid intelligence compared with traditional dispatching rules.

Table 2. Hyperparameters

Hyperparameters	Values	Hyperparameters	Values
Number of training episodes	3000	Discount factor $\gamma$	0.99
Number of testing episodes	500	Prioritized replay $\alpha$	0.6
Schedule cycle	8	Prioritized replay $\beta$	0.4
Buffer size	100000	Number of layers of GRL module	3
Step size	0.0003	Number of attention heads	5
Batch size	64	Random rate of the environment	0.1
Target Q update frequency	100	Shuffle	True

In particular, we show that dynamically selecting appropriate dispatching rules based on the latest problem states can provide consistent superiority over any single fixed dispatching rule. We also compare the proposed method with genetic algorithms (GAs) to demonstrate its performance from a different perspective. Furthermore, thanks to the unified interface provided by Gymjsp, it is more convenient than ever to conduct systematic studies of various RL algorithms on the DJSP. The competitive algorithms adopted in our studies are listed as follows (the same problem formulation, GRL methods, and the number of training episodes are used for all RL methods):

1. Eight dispatching rules: FIFO, LIFO, MOR, LOR, LPT, SPT, LTPT, STPT
2. GA (Chen et al. 2020): GA is a meta-heuristic technique reflecting the principles of biological evolution where the fittest individuals are selected to produce the offspring to form the new generation.
3. A2C (Chen, Hu, and Min 2019): Advantage actor-critic is an asynchronous policy-based RL method, using the actor to interact with the environment and the critic to criticize the actions made by the actor.
4. PPO (Schulman et al. 2017): Proximal policy optimization is a highly effective policy gradient RL method built upon trust region policy optimization (TRPO), which can obtain monotonic improvement by controlling the divergence between the old policy and the new policy.
5. DQN (Luo 2020): Deep Q-Network combines Q-learning with convolutional neural networks and experience replay, well regarded as the foundation of deep RL research.
6. Rainbow DQN (Hessel et al. 2018): Rainbow combines several independent extensions to DQN to provide state-of-the-art performance.

For the sake of generalization and fairness, we keep the efforts on hyperparameter tuning to the minimum. Table 2 summarizes the hyperparameters of D3QPN, whose values are well known to be appropriate on various tasks and are kept unchanged across all instances.

We measure the performance of D3QPN on 12 instances of different sizes to explore the flexibility of D3QPN. All methods are evaluated for 500 episodes, and the full results are shown in Table 3. In general, our method achieves the minimum makespan on all instances except ft06. The reason is that, with **schedule cycle** set to 8, our method can only make 7 or 8 times decisions on ft06, which dramatically limits its performance. If we allow D3QPN to make a decision in every step, it can achieve the same performance as the GA on ft06. Compared with dispatching rules, our proposed method achieves smaller makespan than any individual dispatching rule, with an average performance improvement of 17.80%, confirming the value and superior-

ity of hybrid intelligence over single-agent decision making. Our method also achieves smaller makespan than the meta-heuristic algorithm GA, and the average performance improvement is 28.51%. Furthermore, D3QPN is more effective than other RL algorithms such as A2C, PPO, and DQN, with an average performance improvement of 11.40%. An interesting finding is that D3QPN can outperform Rainbow, which contains more extensions to DQN than D3QPN. In the ablation study in Section 5.5, we will conduct further analysis to show that not all the components in Rainbow are beneficial for the DJSP.

Table 3. The comparison of various algorithms on DJSP instances. The size of the instance is given by the number of jobs  $\times$  the number of machines, and the best makespan results are printed in boldface.

Instance	Size	Dispatching rule								Reinforcement Learning					
		FIFO	LIFO	LPT	SPT	LTPT	STPT	MOR	LOR	GA	A2C	PPO	DQN	Rainbow	Ours
ft06	6 $\times$ 6	64	70	74	85	66	78	59	67	<b>58</b>	69	67	65	63	60
la01	10 $\times$ 5	826	792	818	762	827	933	780	939	738	830	828	785	935	<b>695</b>
la06	15 $\times$ 5	996	1234	1119	1181	1074	1094	927	1097	982	1043	1021	984	1066	<b>921</b>
la11	20 $\times$ 5	1247	1496	1468	1432	1418	1493	1227	1570	1330	1225	1331	1283	1480	<b>1202</b>
la21	15 $\times$ 10	1304	1379	1406	1304	1338	1473	1267	1480	1502	1334	1345	1347	1494	<b>1210</b>
la31	30 $\times$ 10	1884	2197	2248	1981	2089	2297	1833	2196	2436	2075	2047	1958	1846	<b>1775</b>
la36	15 $\times$ 15	1620	1833	1799	1809	1747	1894	1490	1932	1873	1692	1682	1642	1828	<b>1481</b>
orb01	10 $\times$ 10	1374	1355	1396	1389	1299	1427	1314	1385	1432	1344	1343	1327	1473	<b>1141</b>
swv01	20 $\times$ 10	2061	1832	2202	1683	1962	1843	2008	1877	2319	1979	1986	1962	2061	<b>1668</b>
swv06	20 $\times$ 15	2455	2250	2597	2232	2312	2419	2321	2355	2960	2369	2354	2311	2333	<b>2068</b>
swv11	50 $\times$ 10	4448	3771	4624	3626	3870	3825	4602	3957	5354	3819	4103	3998	3937	<b>3534</b>
yn1	20 $\times$ 20	1158	1181	1139	1121	1152	1196	1053	1238	1496	1250	1132	1109	1110	<b>1032</b>

### 5.3. The effect of state representation

To verify the rationality of using attention mechanism as the GRL method, we compare attention mechanism with GNN and matrix representation, and the ranking results are shown in Figure 8. It is clear that the attention mechanism achieves better performance than the GNN and the matrix representation on all instances. To gain a deep insight into its performance, we visualize the attention weights in different dimensions, and the resulting attention maps are shown in Figure 9.

Each map contains the attention weights of each operation (x-axis) relative to all other operations (y-axis). For example, the element in the first row and the second column represents the attention weight that  $O_{00}$  assigns to  $O_{01}$ , which is the amount of attention that  $O_{00}$  should pay to  $O_{01}$  (the darker the square, the greater the weight). According to the visualization results, there is a list of important findings:

- From (a), (c), and (d): the first layer (shallow layer) of the attention model pays more attention to the subsequent operations of each job.
- By comparing (a) and (b): in the first layer, the unfinished operations pay more attention to the subsequent unfinished operations and ignore the completed operations; completed operations pay more attention to completed operations, which is a simple rule to determine the current processing operation of each job.
- By comparing (a) and (e) or (a) and (g) or (b) and (f): unlike the attention in shallow layers where each operation only focuses on the subsequent operations, deeper layers develop a more complex model that can be seen as a kind of more abstract knowledge according to the current operating state.

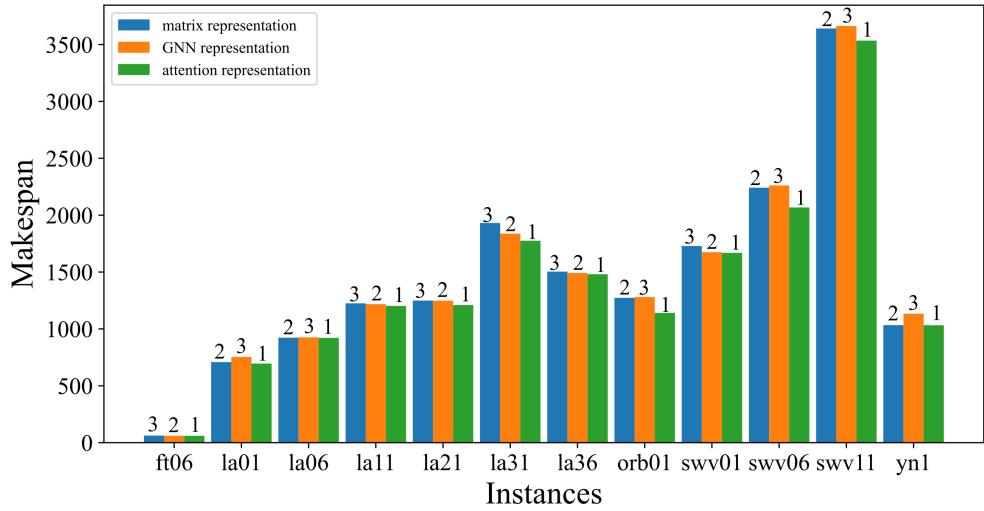


Figure 8. The Makespan performance of three GRL methods over different instances. The numbers above the bar show the ranks on each instance.

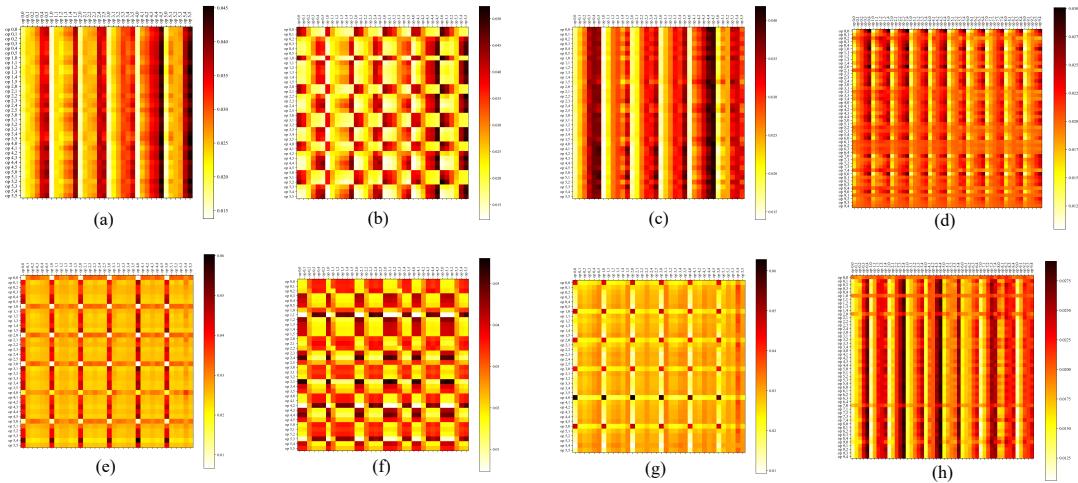


Figure 9. Visualization of attention weights at the convergence of the GRL module. (a) the attention map of the first layer when the input is the initial state of ft06; (b) the attention map of the first layer when the input is the state of ft06 after 20 steps; (c) the attention map of the first layer when the input is the initial state of ft06 with random initialization; (d) the attention map of the first layer when the input is the initial state of la01; (e) the attention map of the third layer when the input is the initial state of ft06; (f) the attention map of the third layer when the input is the state of ft06 after 20 steps; (g) the attention map of another head of the third layer when the input is the initial state of ft06. (h) the attention map of the first layer of the GRL module trained on ft06 when the input is the initial state of la01.

- From (e) and (g): different heads of attention use different strategies to assign weights, confirming that multi-head attention is beneficial for acquiring valuable information from different angles.
- By comparing (a), (c), and (h): it shows that the attention mechanism is robust and can adapt to dynamic events.

#### 5.4. The effect of action space

In some RL-based DJSP methods, customized dispatching rules are used as the action space due to their clarity, ease of implementation, and good interpretability. However, there is a lack of empirical evidence on the advantage over using unrestrained operations or eligible operations as the action space. Figure 10 presents the experimental results on the effect of different action spaces.

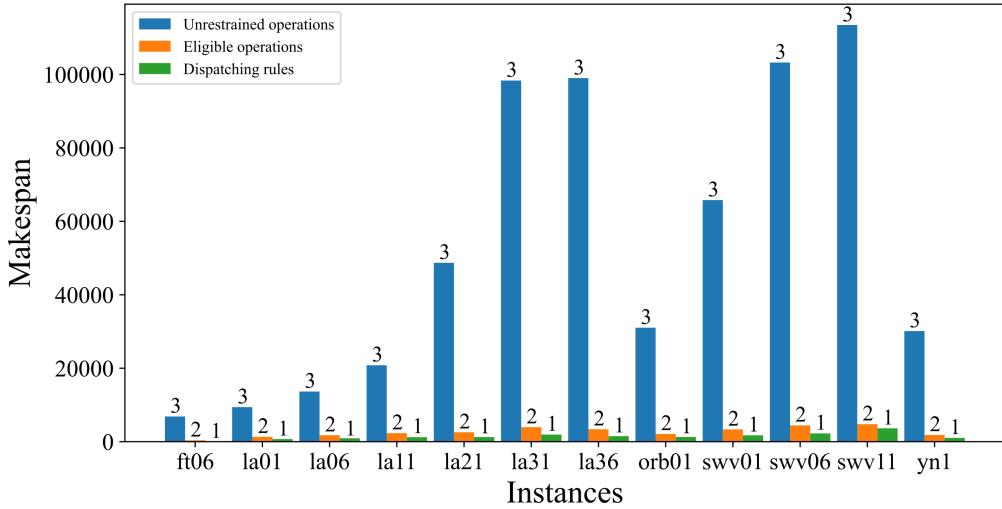


Figure 10. The Makespan performance using three types of action spaces over different instances. The numbers above the bar show the ranks on each instance.

It is clear that, on the 12 instances, using dispatching rules as the actions is always better than using eligible or unrestrained operations. Moreover, using unrestrained operations as actions produce significantly worse results than others. This is because, without any restrictions on the selection of operations, the agent is likely to select a non-executable operation that has been completed, is being processed, requires an occupied machine, or whose antecedent operation has not been completed. Compared with selecting an eligible operation in each step, selecting a dispatching rule to be executed for some time cycles can ensure the consistence of operations and obtain more stable performance. More importantly, using dispatching rules as the action space provides a principled approach to unifying structured domain knowledge and data-driven intelligence with improved flexibility and efficiency.

#### 5.5. Ablation study

As shown in Section 5.2, the performance of D3QPN is significantly better than DQN and Rainbow. To better understand the contribution of each component in D3QPN,

Table 4. The makespan performance of DQN, D3QPN and DQN variants each with single component from Rainbow.  $\uparrow$  and  $\downarrow$  indicate that this method has larger and smaller makespan than DQN, respectively. The results printed in boldface indicate better performance than DQN.

Instance	Makespan-Improvement								D3QPN
	DQN	Double	Dueling	Per	Noisy	Cat	N-step		
ft06	65	<b>64</b> $\downarrow 1\%$	<b>61</b> $\downarrow 6\%$	<b>64</b> $\downarrow 2\%$	<b>59</b> $\downarrow 8\%$	<b>64</b> $\downarrow 1\%$	<b>62</b> $\downarrow 4\%$	<b>60</b> $\downarrow 7\%$	
la01	785	<b>709</b> $\downarrow 10\%$	<b>780</b> $\downarrow 1\%$	<b>783</b> $0\%$	<b>718</b> $\downarrow 9\%$	844 $\uparrow 8\%$	814 $\uparrow 4\%$	<b>695</b> $\downarrow 11\%$	
la06	984	<b>971</b> $\downarrow 1\%$	<b>965</b> $\downarrow 2\%$	996 $\uparrow 1\%$	<b>959</b> $\downarrow 3\%$	1076 $\uparrow 9\%$	<b>969</b> $\downarrow 1\%$	<b>921</b> $\downarrow 6\%$	
la11	1283	<b>1278</b> $0\%$	<b>1234</b> $\downarrow 4\%$	<b>1238</b> $\downarrow 4\%$	<b>1229</b> $\downarrow 4\%$	1299 $\uparrow 1\%$	<b>1240</b> $\downarrow 3\%$	<b>1202</b> $\downarrow 6\%$	
la21	1347	<b>1279</b> $\downarrow 5\%$	<b>1297</b> $\downarrow 4\%$	<b>1306</b> $\downarrow 3\%$	<b>1299</b> $\downarrow 4\%$	<b>1339</b> $\downarrow 1\%$	1351 $0\%$	<b>1210</b> $\downarrow 10\%$	
la31	1958	<b>1928</b> $\downarrow 2\%$	<b>1893</b> $\downarrow 3\%$	<b>1914</b> $\downarrow 2\%$	<b>1920</b> $\downarrow 2\%$	2092 $\uparrow 7\%$	2110 $\uparrow 8\%$	<b>1775</b> $\downarrow 9\%$	
la36	1642	<b>1639</b> $0\%$	1648 $0\%$	<b>1570</b> $\downarrow 4\%$	<b>1628</b> $\downarrow 1\%$	1705 $\uparrow 4\%$	1885 $\uparrow 15\%$	<b>1481</b> $\downarrow 10\%$	
orb01	1327	<b>1267</b> $\downarrow 5\%$	<b>1230</b> $\downarrow 7\%$	<b>1195</b> $\downarrow 10\%$	<b>1161</b> $\downarrow 12\%$	1374 $\uparrow 4\%$	1355 $\uparrow 2\%$	<b>1141</b> $\downarrow 14\%$	
swv01	1962	<b>1772</b> $\downarrow 10\%$	<b>1770</b> $\downarrow 10\%$	<b>1770</b> $\downarrow 10\%$	<b>1711</b> $\downarrow 13\%$	2080 $\uparrow 6\%$	<b>1933</b> $\downarrow 2\%$	<b>1668</b> $\downarrow 15\%$	
swv06	2311	<b>2244</b> $\downarrow 3\%$	<b>2144</b> $\downarrow 7\%$	2337 $\uparrow 1\%$	<b>2299</b> $\downarrow 1\%$	2339 $\uparrow 1\%$	2377 $\uparrow 3\%$	<b>2068</b> $\downarrow 11\%$	
swv11	3998	<b>3881</b> $\downarrow 3\%$	<b>3632</b> $\downarrow 9\%$	<b>3635</b> $\downarrow 9\%$	<b>3697</b> $\downarrow 8\%$	<b>3891</b> $\downarrow 3\%$	<b>3661</b> $\downarrow 8\%$	<b>3534</b> $\downarrow 12\%$	
yn1	1109	<b>1089</b> $\downarrow 2\%$	1047 $\downarrow 6\%$	<b>1055</b> $\downarrow 5\%$	1124 $\uparrow 1\%$	<b>1095</b> $\downarrow 1\%$	1114 $0\%$	<b>1032</b> $\downarrow 7\%$	
Average	1564	<b>1510</b> $\downarrow 3\%$	<b>1475</b> $\downarrow 5\%$	<b>1488</b> $\downarrow 4\%$	<b>1484</b> $\downarrow 5\%$	1600 $\uparrow 3\%$	1573 $\uparrow 1\%$	<b>1399</b> $\downarrow 10\%$	

we perform an ablation study in which we investigate the performance of a number of DQN variants each with a single component in Rainbow. The experiment results are shown in Table 4.

We can see that not all the variants in Rainbow achieve better performance than DQN, which means that some of them are not suitable for the DJSP. For example, distributional DQN (Bellemare, Dabney, and Munos 2017) and multi-step DQN (Peng and Williams 1994) achieve higher makespan values than DQN in general, although multi-step DQN performs better than DQN on some instances. This explains why Rainbow’s performance is worse than that of D3QPN: the multi-step learning and distributional network are negative factors on Rainbow. We hypothesize that this is because the distributional network inherently provides the agent with a distributional perspective on action values, but there are just a few actions to choose from in the DJSP, creating more interference than necessary. Furthermore, the multi-step learning considers the reward before many steps. However, we would prefer the agent to focus on the immediate reward, the machine utilization, rather than the previous reward. In contrast, each component in D3QPN improves the performance of DQN, with NoisyDQN and DuelingDQN, in particular, reducing makespan by 5%.

We also visualize the cumulative rewards during the training of each method to further support our conclusion. The results in Figure 11 show that D3QPN achieves maximal cumulative rewards on almost all instances with a stable decision-making policy. Also, the components in D3QPN help DQN achieve better performance, while the performance of the distributional DQN becomes worse as the training goes on.

## 6. Conclusion

In this paper, we propose a flexible RL-based framework for solving DJSP tasks, which takes disjunctive graphs as states and a set of general dispatching rules as the action space. It is a promising paradigm of hybrid intelligence as it combines structured human knowledge in the form of dispatching rules and flexible data-driven machine intelligence. Moreover, we use the attention mechanism as the graph representation

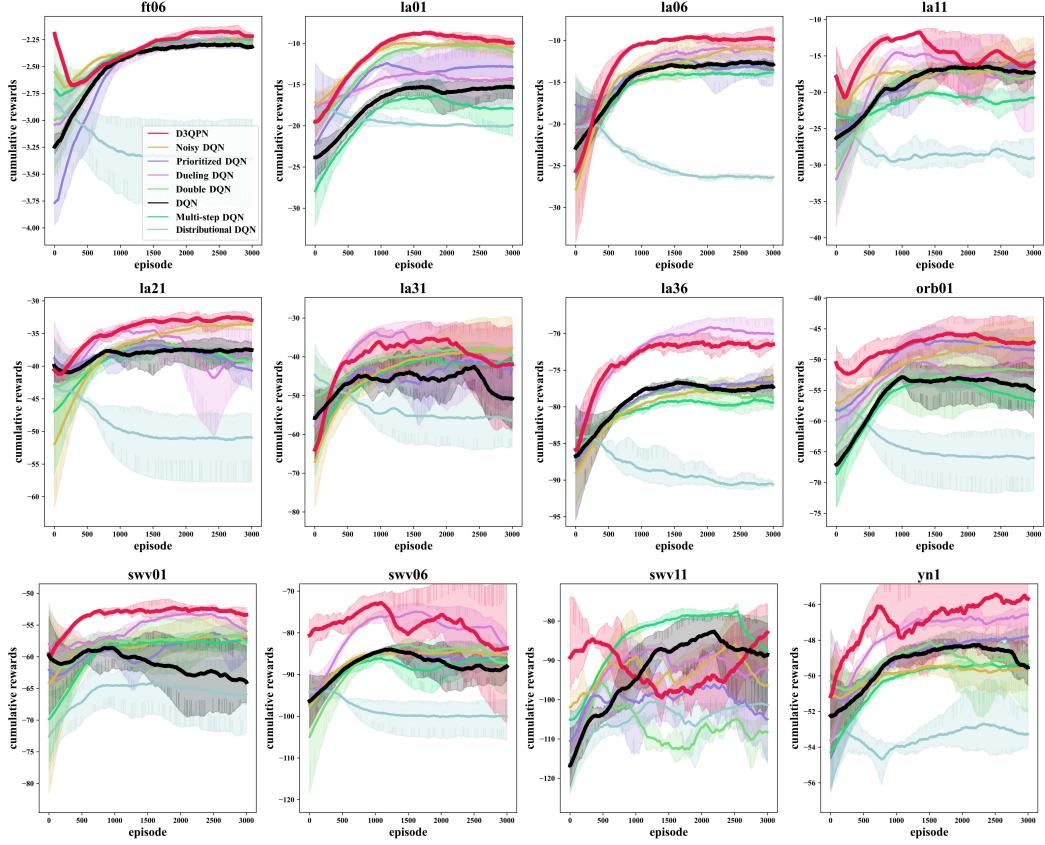


Figure 11. The training curves of DQN, D3QPN, and DQN variants each with a single component from Rainbow. The vertical axis of the graph shows the method’s cumulative reward every episode, reflecting the performance of each method. The solid lines and shaded regions represent the mean and standard deviation, respectively, across five runs.

learning module for effective feature extraction and present D3QPN, an improved DQN algorithm, to map each DJSP state to the most appropriate dispatching rule. To promote the application of RL in the JSP, we present a public benchmark Gymjsp based on the well-known OR-Library, which provides a standardized off-the-shelf and easy-to-use suite for RL communities. Experimental results based on Gymjsp show that our proposed method can outperform traditional rule-based decision systems, meta-heuristic techniques, and SOTA RL algorithms, with great potential in the field of smart manufacturing.

As to future work, it is worthwhile to extend our framework to other variants of JSP, such as the FJSP. Due to the flexibility of our framework, it is also possible to investigate various novel JSP scenarios where the machines or jobs may exhibit unprecedented levels of dynamics. Furthermore, specialized domain knowledge can be strategically integrated into our framework to, for example, design more indicative reward functions or extract more valuable state features.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## Data Availability Statement

The data that support the findings of this study are openly available in <https://github.com/Yunhui1998/Gymjsp>.

## Funding

This work is supported by the National Natural Science Foundation of China (U1713214).

## Notes on contributor(s)

**Yunhui Zeng:** Methodology, Conceptualization, Writing - Original Draft, Coding, Formal analysis. **Zijun Liao:** Coding, Writing - Review & Editing, Formal analysis. **Yuanzhi Dai:** Investigation, Writing - Review & Editing. **Rong Wang:** Investigation, Writing - Review & Editing. **Xiu Li:** Supervision, Writing - Review & Editing. **Bo Yuan.** Resources, Supervision, Writing - Review & Editing.

## References

- Arxiv, Kfir, Helman Stern, and Yael Edan. 2016. “Collaborative reinforcement learning for a two-robot job transfer flow-shop scheduling problem.” *International Journal of Production Research* 54 (4): 1196–1209.
- Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. “Layer normalization.” *arXiv preprint arXiv:1607.06450* .
- Balas, Egon. 1969. “Machine sequencing via disjunctive graphs: an implicit enumeration algorithm.” *Operations Research* 17 (6): 941–957.
- Beasley, John E. 1990. “OR-Library: distributing test problems by electronic mail.” *Journal of the Operational Research Society* 41 (11): 1069–1072.
- Bellemare, Marc G, Will Dabney, and Rémi Munos. 2017. “A distributional perspective on reinforcement learning.” In *International Conference on Machine Learning*, 449–458. PMLR.
- Brockman, Greg, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. “Openai gym.” *arXiv preprint arXiv:1606.01540* .
- Chen, Ronghua, Bo Yang, Shi Li, and Shilong Wang. 2020. “A self-learning genetic algorithm based on reinforcement learning for flexible job-shop scheduling problem.” *Computers & Industrial Engineering* 149: 106778.
- Chen, Xinlei, Li-Jia Li, Li Fei-Fei, and Abhinav Gupta. 2018. “Iterative visual reasoning beyond convolutions.” In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7239–7248.
- Chen, Zheyi, Jia Hu, and Geyong Min. 2019. “Learning-based resource allocation in cloud data center using advantage actor-critic.” In *ICC 2019-2019 IEEE International Conference on Communications (ICC)*, 1–6. IEEE.
- Cui, Zhiyong, Kristian Henrickson, Ruimin Ke, and Yinhai Wang. 2019. “Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic

- learning and forecasting.” *IEEE Transactions on Intelligent Transportation Systems* 21 (11): 4883–4894.
- Das, Sumit, Aritra Dey, Akash Pal, and Nabamita Roy. 2015. “Applications of artificial intelligence in machine learning: Review and prospect.” *International Journal of Computer Applications* 115 (9).
- Demirkol, Ebru, Sanjay Mehta, and Reha Uzsoy. 1998. “Benchmarks for shop scheduling problems.” *European Journal of Operational Research* 109 (1): 137–141.
- Deng, Li, and Dong Yu. 2014. “Deep learning: Methods and applications.” *Foundations and Trends in Signal Processing* 7 (3–4): 197–387.
- Fortunato, Meire, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, et al. 2017. “Noisy networks for exploration.” *arXiv preprint arXiv:1706.10295*.
- Garey, Michael R, David S Johnson, and Ravi Sethi. 1976. “The complexity of flowshop and jobshop scheduling.” *Mathematics of Operations Research* 1 (2): 117–129.
- Gonçalves, José Fernando, Jorge José de Magalhães Mendes, and Mauricio GC Resende. 2005. “A hybrid genetic algorithm for the job shop scheduling problem.” *European Journal of Operational Research* 167 (1): 77–95.
- Hameed, Mohammed Sharafath Abdul, and Andreas Schwung. 2020. “Reinforcement learning on job shop scheduling problems using graph networks.” *arXiv preprint arXiv:2009.03836*.
- Hamilton, William L. 2020. “Graph representation learning.” *Synthesis Lectures on Artificial Intelligence and Machine Learning* 14 (3): 1–159.
- Haupt, Reinhard. 1989. “A survey of priority rule-based scheduling.” *Operations-Research-Spektrum* 11 (1): 3–16.
- He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. “Deep residual learning for image recognition.” In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- Hessel, Matteo, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. 2018. “Rainbow: Combining improvements in deep reinforcement learning.” In *Thirty-second AAAI Conference on Artificial Intelligence*, .
- Kendall, Alex, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John-Mark Allen, Vinh-Dieu Lam, Alex Bewley, and Amar Shah. 2019. “Learning to drive in a day.” In *2019 International Conference on Robotics and Automation (ICRA)*, 8248–8254. IEEE.
- Lin, Chun-Cheng, Der-Jiunn Deng, Yen-Ling Chih, and Hsin-Ting Chiu. 2019. “Smart manufacturing scheduling with edge computing using multiclass deep Q network.” *IEEE Transactions on Industrial Informatics* 15 (7): 4276–4284.
- Liu, Chien-Liang, Chuan-Chin Chang, and Chun-Jan Tseng. 2020. “Actor-critic deep reinforcement learning for solving job shop scheduling problems.” *IEEE Access* 8: 71752–71762.
- Luo, Bin, Sibao Wang, Bo Yang, and Lili Yi. 2021. “An improved deep reinforcement learning approach for the dynamic job shop scheduling problem with random job arrivals.” In *Journal of Physics: Conference Series*, Vol. 1848, 012029. IOP Publishing.
- Luo, Shu. 2020. “Dynamic scheduling for flexible job shop with new job insertions by deep reinforcement learning.” *Applied Soft Computing* 91: 106208.
- Manne, Alan S. 1960. “On the job-shop scheduling problem.” *Operations Research* 8 (2): 219–223.
- Mohan, Jatoth, Krishnanand Lanka, and A Neelakanteswara Rao. 2019. “A review of dynamic job shop scheduling techniques.” *Procedia Manufacturing* 30: 34–39.
- Park, Junyoung, Jaehyeong Chun, Sang Hun Kim, Youngkook Kim, and Jinkyoo Park. 2021. “Learning to schedule job-shop problems: Representation and policy learning using graph neural network and reinforcement learning.” *International Journal of Production Research* 59 (11): 3360–3377.
- Peng, Jing, and Ronald J Williams. 1994. “Incremental multi-step Q-learning.” In *Machine Learning Proceedings 1994*, 226–232. Elsevier.
- Qu, Shuhui, Jie Wang, and Govil Shivani. 2016. “Learning adaptive dispatching rules for a

- manufacturing process system by using reinforcement learning approach.” In *2016 IEEE 21st International Conference on Emerging Technologies and Factory Automation (ETFA)*, 1–8. IEEE.
- Rauber, Paulo, Avinash Ummadisingu, Filipe Mutz, and Jürgen Schmidhuber. 2021. “Reinforcement Learning in Sparse-Reward Environments With Hindsight Policy Gradients.” *Neural Computation* 33 (6): 1498–1553.
- Reyna, Yunior César Fonseca, and Yailen Martínez-Jiménez. 2017. “Adapting a Reinforcement Learning Approach for the Flow Shop Environment with sequence-dependent setup time.” *Revista Cubana de Ciencias Informáticas* 11 (1): 41–57.
- Richter, Florian, Ryan K Orosco, and Michael C Yip. 2019. “Open-sourced reinforcement learning environments for surgical robotics.” *arXiv preprint arXiv:1903.02090*.
- Satyro, Walter Cardoso, Mauro de Mesquita Spinola, Cecília MVB de Almeida, Biagio F Giannetti, José Benedito Sacomano, Jose Celso Contador, and Jose Luiz Contador. 2021. “Sustainable industries: Production planning and control as an ally to implement strategy.” *Journal of Cleaner Production* 281: 124781.
- Schaul, Tom, John Quan, Ioannis Antonoglou, and David Silver. 2015. “Prioritized experience replay.” *arXiv preprint arXiv:1511.05952*.
- Schulman, John, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. “Proximal policy optimization algorithms.” *arXiv preprint arXiv:1707.06347*.
- Senior, Anderw, John Jumper, Demis Hassabis, and Pushmeet Kohli. 2018. “AlphaFold: Using AI for scientific discovery.” *DeepMind. Recuperado de: https://deepmind. com/blog/alphafold*
- Shi, Min, David A Wilson, Xingquan Zhu, Yu Huang, Yuan Zhuang, Jianxun Liu, and Yufei Tang. 2020. “Evolutionary architecture search for graph neural networks.” *arXiv preprint arXiv:2009.10199*.
- Silver, David, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, et al. 2016. “Mastering the game of Go with deep neural networks and tree search.” *Nature* 529 (7587): 484–489.
- Silver, David, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, et al. 2017. “Mastering the game of go without human knowledge.” *Nature* 550 (7676): 354–359.
- Sutton, Richard S, and Andrew G Barto. 2018. *Reinforcement learning: An introduction*. MIT press.
- Taillard, Eric. 1993. “Benchmarks for basic scheduling problems.” *European Journal of Operational Research* 64 (2): 278–285.
- Teymourifar, Aydin, Gurkan Ozturk, Zehra Kamisli Ozturk, and Ozan Bahadir. 2020. “Extracting new dispatching rules for multi-objective dynamic flexible job shop scheduling with limited buffer spaces.” *Cognitive Computation* 12 (1): 195–205.
- Tian, W, and HP Zhang. 2021. “A dynamic job-shop scheduling model based on deep learning.” *Advances in Production Engineering & Management* 16 (1).
- Turgut, Yakup, and Cafer Erhan Bozdag. 2020. “Deep Q-network model for dynamic job shop scheduling pproblem based on discrete event simulation.” In *2020 Winter Simulation Conference (WSC)*, 1551–1559. IEEE.
- Van Hasselt, Hado, Arthur Guez, and David Silver. 2016. “Deep reinforcement learning with double q-learning.” In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 30.
- Van Laarhoven, Peter JM, Emile HL Aarts, and Jan Karel Lenstra. 1992. “Job shop scheduling by simulated annealing.” *Operations Research* 40 (1): 113–125.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. “Attention is all you need.” In *Advances in Neural Information Processing Systems*, 5998–6008.
- Wang, Haoxiang, Bhaba R Sarker, Jing Li, and Jian Li. 2020. “Adaptive scheduling for assembly job shop with uncertain assembly times based on dual Q-learning.” *International Journal of Production Research* 1–17.
- Wang, Libing, Xin Hu, Yin Wang, Sujie Xu, Shijun Ma, Kexin Yang, Zhijun Liu, and Weidong

- Wang. 2021. “Dynamic job-shop scheduling in smart manufacturing using deep reinforcement learning.” *Computer Networks* 190: 107969.
- Wang, Yi-Chi, and John M Usher. 2005. “Application of reinforcement learning for agent-based production scheduling.” *Engineering Applications of Artificial Intelligence* 18 (1): 73–82.
- Wang, Ziyu, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. 2016. “Dueling network architectures for deep reinforcement learning.” In *International conference on machine learning*, 1995–2003. PMLR.
- Wei, Y, and M Zhao. 2005. “A reinforcement learning-based approach to dynamic job-shop scheduling.” *Acta Automatica Sinica* 31 (5): 765.
- Wu, Zonghan, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. 2020. “A comprehensive survey on graph neural networks.” *IEEE Transactions on Neural Networks and Learning Systems* 32 (1): 4–24.
- Xia, Kaishu, Christopher Sacco, Max Kirkpatrick, Clint Saidy, Lam Nguyen, Anil Kircaliali, and Ramy Harik. 2021. “A digital twin to train deep reinforcement learning agent for smart manufacturing plants: Environment, interfaces and intelligence.” *Journal of Manufacturing Systems* 58: 210–230.
- Yang, H-B, and H-S Yan. 2009. “An adaptive approach to dynamic scheduling in knowledgeable manufacturing cell.” *The International Journal of Advanced Manufacturing Technology* 42 (3): 312–320.
- Zhang, Cong, Wen Song, Zhiguang Cao, Jie Zhang, Puay Siew Tan, and Chi Xu. 2020. “Learning to dispatch for job shop scheduling via deep reinforcement learning.” *arXiv preprint arXiv:2010.12367* .
- Zhao, Y, and H Zhang. 2021. “Application of machine learning and rule scheduling in a job-shop production control system.” *International Journal of Simulation Modelling (IJSIMM)* 20 (2).
- Zhou, Jie, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2020. “Graph neural networks: A review of methods and applications.” *AI Open* 1: 57–81.