

Deep Reinforcement Learning for Job Scheduling and Resource Management in Cloud Computing: An Algorithm-Level Review

Yan Gu, Zhaoze Liu, Shuhong Dai, Cong Liu, Ying Wang, Shen Wang, Georgios Theodoropoulos, Long Cheng

Abstract—Cloud computing has revolutionized the provisioning of computing resources, offering scalable, flexible, and on-demand services to meet the diverse requirements of modern applications. At the heart of efficient cloud operations are job scheduling and resource management, which are critical for optimizing system performance and ensuring timely and cost-effective service delivery. However, the dynamic and heterogeneous nature of cloud environments presents significant challenges for these tasks, as workloads and resource availability can fluctuate unpredictably. Traditional approaches, including heuristic and meta-heuristic algorithms, often struggle to adapt to these real-time changes due to their reliance on static models or pre-defined rules. Deep Reinforcement Learning (DRL) has emerged as a promising solution to these challenges by enabling systems to learn and adapt policies based on continuous observations of the environment, facilitating intelligent and responsive decision-making. This survey provides a comprehensive review of DRL-based algorithms for job scheduling and resource management in cloud computing, analyzing their methodologies, performance metrics, and practical applications. We also highlight emerging trends and future research directions, offering valuable insights into leveraging DRL to advance both job scheduling and resource management in cloud computing.

Index Terms—Deep reinforcement learning, cloud computing, job scheduling, resource management, survey

I. INTRODUCTION

Cloud Computing. Cloud computing has fundamentally reshaped the landscape of modern computing, offering flexible, scalable, and cost-effective solutions for data storage, processing, and management. Unlike traditional computing models, where users rely on local servers or on-premises infrastructure, cloud computing provides a distributed environment where computational resources, including servers, storage, and software, are delivered over the internet [1]. This shift to cloud-based services enables organizations and individuals to access powerful computing resources without the need to invest heavily in physical infrastructure, allowing them to scale their operations up or down based on demand.

Y. Gu, Z. Liu, S. Dai and L. Cheng are with the School of Control and Computer Engineering, North China Electric Power University in Beijing, China. E-mail: lcheng@ncepu.edu.cn

C. Liu is with the School of Computer Science and Technology, Shandong University of Technology, China. E-mail: liucongchina@sdust.edu.cn

Y. Wang is with the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China. E-mail: wangying2009@ict.ac.cn

S. Wang is with the School of Computer Science, University College Dublin, Ireland, E-mail: shen.wang@ucd.ie

G. Theodoropoulos is with Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen, China. E-mail: georgios@sustc.edu.cn

Cloud computing is underpinned by a variety of service models, including Infrastructure-as-a-Service (IaaS), Platform-as-a-Service (PaaS), and Software-as-a-Service (SaaS), each offering different levels of abstraction and control. These services support a wide array of applications, from enterprise resource planning (ERP) systems to machine learning platforms [2], [3], and have become critical enablers of innovation across industries such as finance, healthcare, and entertainment. As the demand for cloud services continues to grow, cloud providers face the ongoing challenge of managing an increasingly complex and dynamic environment to ensure high performance, reliability, and efficiency.

As cloud computing expands to incorporate edge computing [4], and scales to support an ever-growing number of users and applications, effective job scheduling and resource management become critical for ensuring optimal performance and resource utilization. Job scheduling involves the allocation of tasks or jobs to available computing resources in a manner that maximizes efficiency, minimizes response times, and ensures fairness among users. Resource management, on the other hand, focuses on the allocation and optimization of computational resources such as CPUs, memory, storage, and bandwidth to meet the diverse needs of various applications and workloads [5].

Job Scheduling and Resource Management. The dynamic and heterogeneous nature of cloud environments makes job scheduling and resource management particularly challenging. Cloud workloads vary significantly in terms of computational intensity, real-time constraints, and data dependencies, which complicates the scheduling process [6], [7]. Furthermore, resources are often distributed across multiple servers, data centers, and geographic locations, making it difficult to ensure consistent performance and effective load balancing. The demand for resources fluctuates based on factors such as workload characteristics, user behavior, and external conditions, requiring adaptive management mechanisms [8], [9], [10]. In addition, cloud providers must address critical issues like energy efficiency and fault tolerance [11], [12]. As a result, efficient job scheduling and resource management become crucial not only for maximizing resource utilization but also for minimizing operational costs, improving quality of service (QoS), and ensuring compliance with service level agreements (SLAs). These factors highlight the need for intelligent, adaptive solutions that can handle the inherent complexities of cloud environments.

A multitude of methods have been developed to address

TABLE I: Comparison of existing reviews and surveys on job scheduling and resource management

Reference	Task Scheduling	Workflow Scheduling	Resource Provisioning	Resource Scheduling	Algorithm-Level	Reviewed Method
[13]	-	-	✓	✓	-	Metaheuristic
[14]	-	-	✓	✓	-	Heuristic and DRL
[5]	-	-	✓	✓	-	Heuristic and DRL
[15]	✓	✓	-	-	-	Metaheuristic
[16]	✓	✓	-	-	-	Metaheuristic
[17]	-	-	✓	✓	-	Metaheuristic
[18]	✓	✓	-	✓	-	Metaheuristic
[19]	✓	-	✓	✓	-	Heuristic and DRL
[20]	✓	-	✓	✓	-	DRL
[21]	✓	✓	✓	✓	-	DRL
This survey	✓	✓	✓	✓	✓	DRL

job scheduling and resource management in cloud computing environments. For job scheduling, conventional rule-based policies—such as round-robin [22]—and algorithms like Min-Min [23] have been widely utilized due to their simplicity and ease of implementation. To enhance scheduling efficiency, heuristic algorithms (e.g., heterogeneous earliest finish time [24]) and meta-heuristic algorithms (e.g., genetic algorithms [25], whale optimization algorithm [26]) have been explored for their ability to find near-optimal solutions in complex scheduling scenarios. Similar approaches have been applied to resource management, including heuristic algorithms based on bin packing [27] and Petri nets [28], as well as meta-heuristic methods like genetic algorithms [29] and particle swarm optimization [30], which have demonstrated potential in optimizing resource allocation strategies. Despite their effectiveness, these heuristic and meta-heuristic algorithms exhibit significant limitations in real-time, dynamic cloud environments. Their reliance on prior knowledge and static optimization models renders them less adaptable to the unpredictable nature of cloud systems, especially when task arrival times and resource demands fluctuate rapidly.

To overcome these challenges, Deep Reinforcement Learning (DRL) has emerged as a robust and adaptive alternative for both job scheduling and resource management. DRL combines reinforcement learning with deep neural networks, enabling systems to learn optimal policies through continuous interaction with the environment. By leveraging historical data and real-time feedback, DRL algorithms can make informed decisions based on the current state of the cloud environment, effectively adapting to changes and uncertainties [31]. This adaptive learning process allows DRL to address the complexities inherent in cloud computing, such as dynamic workloads, resource heterogeneity, and unpredictable demand patterns. Consequently, DRL facilitates the development of intelligent scheduling and resource allocation strategies that optimize resource utilization, enhance system performance, and improve QoS.

Related Reviews. In recent years, significant advances in job scheduling and resource management have driven extensive applications in the domain of cloud computing. Building on these advances, many reviews and surveys have summarized various developments in this field. To provide a more comprehensive understanding, we have synthesized the insights from these reviews to present a holistic perspective. As illustrated in Table I, our reviews offers a superior con-

tribution compared to existing works by providing a more comprehensive and detailed analysis of DRL techniques aimed at addressing the challenges of job scheduling and resource management. Specifically, the works [13], [14], [5] focus on resource management in network function virtualization, 5G networks, and edge computing, respectively, while neglecting a thorough investigation of job scheduling. In comparison, the surveys in [15], [16] provide detailed analyses of task and workflow scheduling in cloud environments relying on meta-heuristic algorithms, but they fall short by overlooking resource management considerations. Studies by [17], [18] offer an extensive overview of resource management and task scheduling, however they do not consider DRL-based approaches. Additionally, works [19], [20], [21] review DRL-based methods for resource management and task scheduling in cloud computing. Despite their contributions, these works fall short of providing an algorithm-level review of DRL methods, hindering a thorough understanding of DRL advancements in task scheduling and resource management. Our work bridges this gap by systematically analyzing DRL in job scheduling and resource management from an algorithm-level perspective, providing a structured examination of advancements and methodologies.

Our Contributions. In this review, we present a comprehensive analysis of the applications of DRL in job scheduling and resource management, emphasizing the categorization and design principles of various DRL methods. Furthermore, as many studies in edge-cloud computing and edge computing share similar settings with cloud environments, their methods are often applicable to cloud computing scenarios. Consequently, unless explicitly specified, we broaden the scope of this review to include such works. The structure of this review is illustrated in Fig. 1, and our main contributions are summarized as follows:

- *Algorithm-Level Review and Analysis:* This review provides an in-depth analysis of DRL algorithms, focusing on their applications in job scheduling and resource management within cloud computing. We categorize these algorithms into four main types: value-based methods, policy-based methods, multi-agent approaches, and advanced DRL techniques. This classification not only highlights the unique strengths and capabilities of each approach but also provides a structured framework to tackle the complex challenges inherent in job scheduling and resource management.

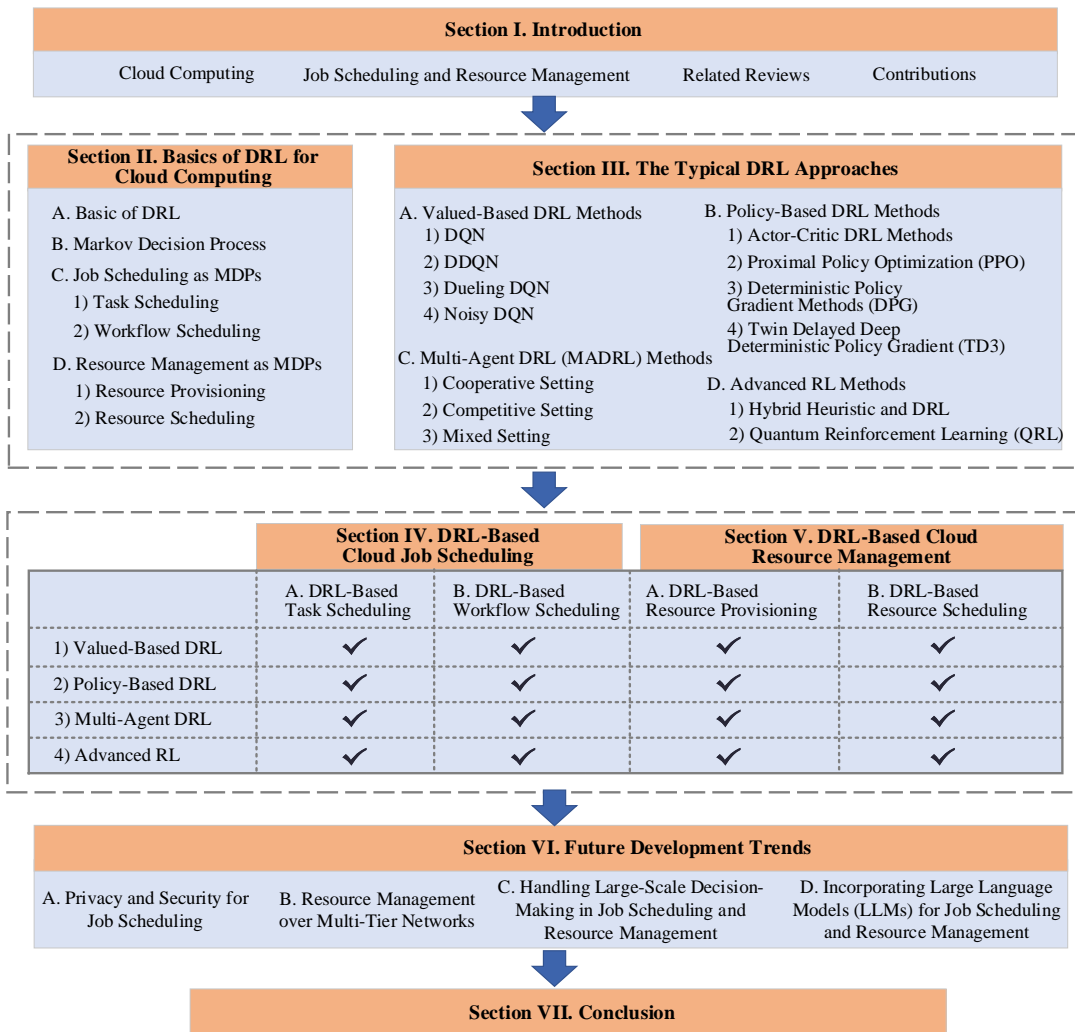


Fig. 1: The general architecture of this review

- *Comprehensive Survey of DRL-based Approaches:* This work provides a detailed review of DRL approaches applied to job scheduling and resource management, including task and workflow scheduling, as well as resource provisioning and allocation strategies. We analyze how DRL techniques are customized to optimize these processes, offering insights into their performance, scalability, and adaptability. Furthermore, we examine the application of DRL in cloud and edge computing environments, highlighting the distinct challenges, opportunities, and considerations specific to each context.
- *Insights into Future Development Trends:* This review highlights critical challenges and opportunities for advancing deep reinforcement learning in job scheduling and resource management within cloud computing. It outlines potential future directions, such as strengthening privacy and security measures, enhancing the robustness and scalability of DRL frameworks, expanding their applicability to dynamic and heterogeneous environments, and improving model interpretability to support practical implementation and real-world adoption.

The remainder of this paper is structured as follows: Sec-

tion II provides an introduction to the fundamentals of deep reinforcement learning and its relevance to cloud computing. Section III offers a comprehensive overview of typical DRL algorithms. Section IV reviews DRL-based methodologies for job scheduling, focusing on both task and workflow scheduling. Section V examines resource management techniques utilizing DRL, including resource provisioning and scheduling strategies. Section VI discusses future directions for advancing DRL in job scheduling and resource management. Finally, Section VII provides the conclusion of this review.

II. BASICS OF DRL FOR CLOUD COMPUTING

In this section, we introduce the fundamentals of DRL and Markov Decision Processes (MDPs), and provide a general introduction to modeling job scheduling and resource management as MDPs, laying the groundwork for applying DRL in cloud computing.

A. Basics of DRL

Reinforcement Learning (RL) has emerged as a powerful paradigm for solving sequential decision-making problems,

where an agent learns optimal behaviors through interactions with an environment [32]. In the RL framework, the agent observes the current state of the environment, selects an action and subsequently receives feedback in the form of a reward which guides future actions. The goal of RL is to develop a policy that maximizes the cumulative reward over time. However, traditional RL methods struggle with high-dimensional state spaces and complex decision problems, especially when state representations and decision-making processes are not straightforward [33].

To address these challenges, DRL combines the principles of RL with the powerful function approximation capabilities of deep neural networks [34]. By using neural networks to approximate value functions or directly model policies, DRL can address problems involving large state and action spaces, commonly found in real-world applications like robotics [35], transportation systems [36], network control [37], and also cloud and edge computing [38].

B. Markov Decision Process

In RL, a learning agent interacts with an environment to address sequential decision-making problems. Fully observable environments are typically modeled as MDPs, which are formally defined by a quintuple $(\mathcal{S}, \mathcal{A}, T, \mathcal{R}, \gamma)$. Here, \mathcal{S} denotes the state space, encompassing all possible states the system can occupy, while \mathcal{A} represents the action space, containing all feasible actions an agent can take in any given state. The transition probability function T defines the likelihood of transitioning from one state to another, expressed as $P(s_{t+1}|s_t, a_t)$, given an action a_t . The reward function \mathcal{R} assigns a scalar reward $R(s_t, a_t)$ based on the agent's action in a specific state, reflecting the immediate value of that action. Finally, the discount factor γ , where $0 \leq \gamma \leq 1$, governs the trade-off between immediate and future rewards, with a higher γ emphasizing long-term gains and a lower γ focusing on short-term outcomes.

At each discrete time step t , the agent observes the current state s_t , selects an action a_t according to a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$, and receives a reward $r_t = R(s_t, a_t)$. The environment then transitions to a new state s_{t+1} based on the transition probability $P(s_{t+1}|s_t, a_t)$. The goal of agent is to identify an optimal policy π^* that maximizes the expected cumulative reward over time, mathematically expressed as:

$$\sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t) \quad (1)$$

where γ determines the relative importance of future rewards in decision-making. The MDP framework provides a foundation for DRL, allowing agents to iteratively improve their policies $\pi(a|s)$ by maximizing the expected cumulative reward, leading to better decision-making over time.

C. Job Scheduling as MDPs

Job scheduling can generally be abstracted into two levels based on the structure of jobs: task scheduling and workflow scheduling. Workflow scheduling can be considered an

extension of task scheduling, as it involves managing more complex dependencies between tasks [39]. In dynamic environments, such as those with rapidly fluctuating job arrival rates, scheduling decisions depend solely on the current system state, with future states determined by immediate scheduling actions. This property makes these scenarios particularly well-suited for modeling as MDPs.

1) *Task Scheduling*: Task scheduling in cloud computing focuses on the allocation of individual tasks to available resources, ensuring that tasks are executed efficiently.

Action Space: Given the sequential arrival of individual tasks over time, the action space at each Markov decision step is designed to schedule the tasks that have arrived at the current moment. Specifically, the scheduler selects an action that determines how the newly arrived task is assigned to available resources or queued for later execution [40]. The action space \mathcal{A} can be mathematically defined as:

$$\mathcal{A} = \{a \mid a = \text{assign } t \text{ to } v, v \in V\} \quad (2)$$

Here, t represents the task to be scheduled at the current decision step, v is a computational resource to which the task can be assigned, and V denotes the set of all available computational resources in the system. In some scenarios, the action space \mathcal{A} may vary depending on task-specific constraints. For instance, due to privacy and security considerations, certain tasks may only be assigned to resources in a private cloud [41]. This dynamic and constraint-aware definition of the action space ensures that scheduling decisions remain feasible and aligned with the specific requirements of tasks and resources.

State Space: The task action space is inherently influenced by the state space \mathcal{S} of the cloud environment due to the Markov property. Broadly, \mathcal{S} comprises two main components: task status S_t and the states of computational resources S_v . This can be expressed as:

$$\mathcal{S} = \{S_t, S_v\} \quad (3)$$

The task status S_t captures critical details about each task, such as CPU, memory, and storage requirements, as well as its execution status (e.g., time remaining for completion). In certain computational scenarios, S_t can be extended to include additional constraints, such as security and privacy requirements or geographical location preferences [42], [43]. For computational resource states, S_v represents the status of available resources, including current availability, expected completion times, computation and storage costs, and other relevant metrics [44].

Reward Function: In cloud systems, task scheduling often involves multiple optimization objectives. The overall reward function is typically defined as: $\mathcal{R} = \mathcal{F}(o_1, o_2, \dots)$, where o_i represents an individual optimization objective. Maximization objectives commonly include factors such as memory utilization, storage utilization, and network bandwidth usage. Conversely, minimization objectives focus on metrics like operational costs, task response time, and energy consumption [45], [46]. In practice, the reward function is often designed as a weighted combination of these objectives, enabling a balanced trade-off between competing goals [47].

2) *Workflow Scheduling*: Cloud workflows are commonly modeled as Directed Acyclic Graphs (DAGs) and are typically decomposed into subworkflows or subtasks [48] during the scheduling process. This decomposition facilitates the parallel execution of workflows across diverse resources in the cloud environment, thereby enhancing resource utilization and workflow execution efficiency.

Action Space: Similar to task scheduling, the action space \mathcal{A} in workflow scheduling can be abstracted at a high level as:

$$\mathcal{A} = \{a \mid a = \text{assign } w \text{ to } v, v \in V\} \quad (4)$$

where w represents the workflow to be scheduled at the current decision step. Generally, the action space in workflow scheduling is typically structured into two decision levels. The first level involves decomposing the workflow w into subworkflows or tasks t , which can be performed using DRL or other ruled-based algorithms [49]. The second level focuses on allocating the decomposed subworkflows or tasks to appropriate resources, taking into account task dependencies within the workflow.

State Space: Unlike task scheduling, which primarily focuses on individual tasks, workflow scheduling involves managing entire workflows efficiently by coordinating the execution of decomposed tasks. As a result, the state space for workflow scheduling is inherently more complex and extends beyond that of task scheduling. The state space \mathcal{S} for workflow scheduling can be represented as:

$$\mathcal{S} = \{S_t, S_w, S_v\} \quad (5)$$

where S_t represents the critical details of the decomposed subworkflows or tasks, and S_w encapsulates the state of the workflows.

Reward Function: In workflow scheduling, the reward function is tied to various optimization objectives, with makespan and cost often being primary considerations. Unlike traditional task scheduling, workflow scheduling introduces unique challenges, particularly in managing costs. In addition to computation and storage costs, communication costs play a critical role due to the data dependencies between tasks executed across distributed resources. These communication costs, which include data transfer time and bandwidth utilization, can significantly influence workflow performance, especially in geographically dispersed environments [50].

D. Resource Management as MDPs

Resource management in cloud computing encompasses two primary functions: resource provisioning and resource scheduling. Resource provisioning involves allocating virtualized resources to accommodate varying workloads and user demands, while resource scheduling focuses on assigning these allocated resources to specific tasks or applications. In dynamic and uncertain cloud environments, resource management decisions are guided by the current state of the system. The evolution of future states depends directly on the actions taken in the present, without reliance on past states, making resource management suitable for representation using MDPs.

1) *Resource Provisioning*: Through elastic scaling of computational resources, resource provisioning enables cloud systems to adapt to fluctuating workload demands, thereby enhancing overall performance and efficiency.

Action Space: In resource provisioning, the action space depends on the type of scaling. For horizontal scaling, the resource v typically refers to virtual machines or containers, while for vertical scaling, the resource v usually represents CPU, memory, or storage capacity [51]. Generally, the action space \mathcal{A} be expressed as:

$$\mathcal{A} = \{a \mid \text{scale up/down/keep unchanged for } v, v \in V\} \quad (6)$$

The actions are often represented as discrete integer values, indicating the number of resources to be scaled up or down.

State Space: To model the resource provisioning problem, the state space \mathcal{S} can be generally expressed as:

$$\mathcal{S} = \{S_u, S_n, S_p\} \quad (7)$$

where S_u represents the current utilization of system resources, such as CPU, memory, storage, and network bandwidth. S_n denotes the number of virtual machine or container instances running on each server. Additionally, S_p captures performance metrics of the system, including throughput, response time, and energy consumption.

Reward Function: The reward function in resource provisioning incorporates multiple optimization objectives include maximizing resource utilization and system throughput, minimizing infrastructure costs and energy consumption, and maintaining load balance to prevent single-point overloads. Each metric is carefully quantified and integrated into the reward function, ensuring a balanced trade-off between performance, cost efficiency, and system stability.

2) *Resource Scheduling*: Resource scheduling aims to allocate resources to jobs efficiently, optimizing objectives such as maximizing resource utilization and system throughput.

Action Space: Give the available resource V and an incoming task t requiring resources at the current step, the action space \mathcal{A} can be defined as:

$$\mathcal{A} = \{a \mid a = \text{allocate } v \text{ to } t, v \in V\} \quad (8)$$

The type of resources varies depending on the specific computing scenario. In cloud environments, available resources typically include virtual machines [52], whereas in edge cloud computing scenarios, the resources could consist of edge devices such as mobile phones [53] and vehicles [54]. Additionally, resources can be subdivided into more granular components, such as processing units (e.g., CPUs, GPUs), storage, memory, and network bandwidth [55], [56].

State Space: Similar to task scheduling, resource scheduling involves both the state of the resources S_v and the state of the tasks S_t . The state space \mathcal{S} can be represented as:

$$\mathcal{S} = \{S_v, S_t\} \quad (9)$$

where S_v captures the current status of resources, and S_t represents the characteristics and requirements of the tasks in the system.

Reward Function: The design of the reward function for resource scheduling typically considers task requirements and

overall system performance. From the task perspective, the primary objectives are to meet QoS and SLA requirements while minimizing the makespan. On the system side, the focus shifts to maximizing resource utilization and minimizing energy consumption and operational costs.

III. THE TYPICAL DRL APPROACHES

In this section, as illustrated in Fig. 2, we provide a detailed overview of typical DRL algorithms applied in cloud computing, establishing a foundation for our algorithm-level review on their use in job scheduling and resource management.

A. Value-Based DRL Methods

Value-based DRL methods center around learning a value function which estimates the expected cumulative reward of a state-action pair. As shown in Fig. 2 (a), these methods typically aim to optimize a policy indirectly by approximating the optimal action-value function $Q^*(s, a)$, which predicts the total expected reward starting from state s , taking action a , and following an optimal policy thereafter. The most widely used value-based method is Deep Q-Network (DQN), which builds upon traditional Q-learning by employing deep neural networks as function approximators [57]. Several variations of DQN have been proposed to address its limitations, including Double DQN (DDQN), Dueling DQN, and Noisy DQN.

1) *DQN*: DQN extends Q-learning by utilizing a deep neural network to approximate the Q-function $Q(s, a; \theta)$ where θ represents the parameters of the network. The agent interacts with the environment, storing experiences (s_t, a_t, r_t, s_{t+1}) in a replay buffer. To update the network, mini-batches of these experiences are sampled and the Q-network is trained to minimize the loss:

$$L(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1})} \left[(y_t - Q(s_t, a_t; \theta))^2 \right] \quad (10)$$

where the target y_t is calculated as:

$$y_t = r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta^-) \quad (11)$$

and θ^- refers to the parameters of a target network updated periodically for stability. At each step, the agent selects the action $a_t = \arg \max_a Q(s_t, a; \theta)$ which maximizes the predicted Q-value for the current state. However, since the same network is used both to select and evaluate actions, this can result in overly optimistic estimates of action values and ultimately lead to overestimation.

2) *DDQN*: DDQN addresses the overestimation issue in DQN by decoupling the action selection and evaluation processes. In this method, the online Q-network selects the action, while the target network evaluates it to reduce overestimation bias [58]. The target for the loss function is updated as:

$$y_t^{\text{Double}} = r_t + \gamma Q(s_{t+1}, \arg \max_{a'} Q(s_{t+1}, a'; \theta); \theta^-) \quad (12)$$

where the online network selects $\arg \max_{a'} Q(s_{t+1}, a'; \theta)$, and the target network θ^- evaluates the Q-value. This decoupling leads to more accurate value estimates and better overall policy learning, particularly in noisy environments.

3) *Dueling DQN*: In many environments, estimating the value of each action can be inefficient. To address this, the dueling architecture decomposes the Q-value into the state value $V(s)$ and the action advantage $A(s, a)$. The Q-value is then expressed as:

$$Q(s, a) = V(s) + A(s, a) \quad (13)$$

To ensure the Q-values are uniquely determined, the mean advantage across all actions is subtracted:

$$Q(s, a) = V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right) \quad (14)$$

This structure allows the agent to focus on learning state values, which improves efficiency in environments where the choice of action has less impact. By decoupling state and action evaluation, Dueling DQN enables faster and more stable learning particularly in large action spaces.

4) *Noisy DQN*: Efficient exploration is essential in DRL, but traditional approaches like ϵ -greedy exploration often fail in environments with complex or sparse reward structures. Noisy DQN introduces learnable noise directly into the parameters of the Q-networks to enhance exploration, eliminating the need for manually tuned exploration schedules and enabling more effective exploration throughout training.

In Noisy DQN, the network weights W are perturbed by adding parameterized noise expressed as:

$$W = \mu + \sigma \cdot \epsilon \quad (15)$$

where μ and σ are learnable parameters and ϵ is drawn from a noise distribution such as Gaussian or factorized noise. This approach ensures that the Q-value estimates vary due to stochasticity in the network and promotes exploration throughout training.

Noisy DQN modifies the standard DQN loss function as described in Equation 10, by introducing noisy weights that enhance exploration through added variability in Q-value estimates without the need for manual tuning. As training progresses, the noise adjusts automatically and enables a smooth shift from exploration to exploitation. This approach has shown better performance than standard DQN particularly in tasks with sparse rewards or large action spaces due to its sustained exploration.

B. Policy-Based DRL Methods

As illustrated in Fig. 2 (b), policy-based DRL methods focus on directly learning a policy $\pi(a|s)$ which specifies the probability of taking action a given state s . Unlike value-based methods that derive a policy indirectly by learning value functions, policy-based methods optimize the policy itself by maximizing the expected cumulative reward over time [59]. These methods are particularly effective in continuous or high-dimensional action spaces, where discretizing the action space for value-based approaches becomes computationally infeasible.

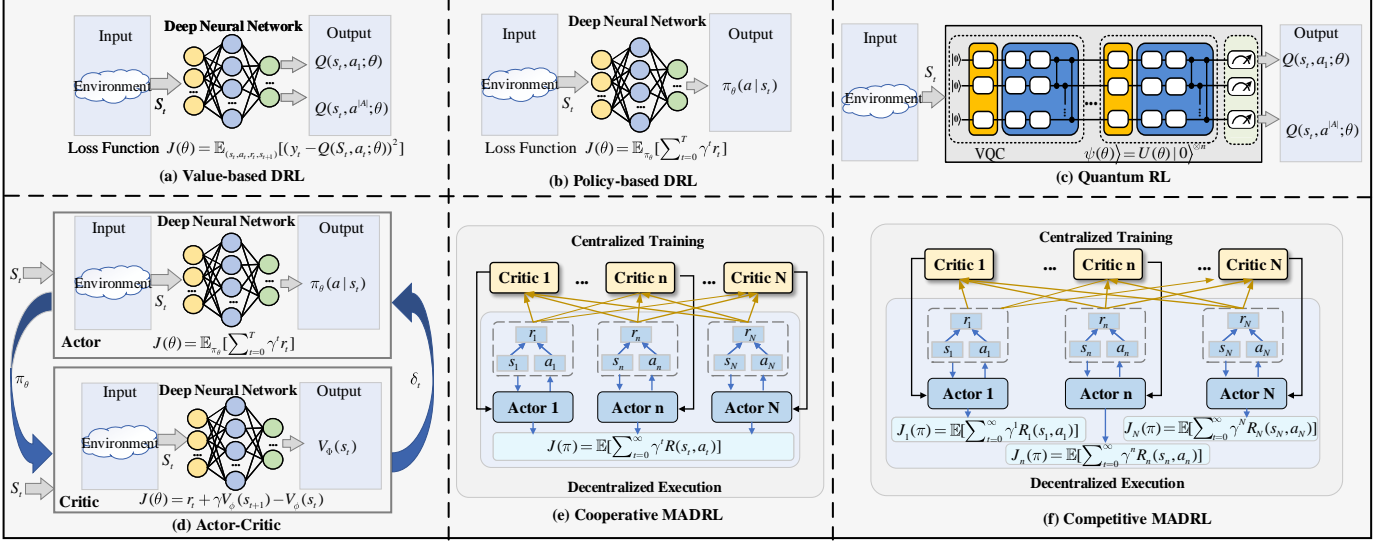


Fig. 2: Detailed architectures of typical DRL approaches

1) *Actor-Critic DRL Methods*: One of the most widely adopted frameworks within policy-based methods is the Actor-Critic (AC) architecture [60]. In this approach, the actor is responsible for learning the policy $\pi_\theta(a|s)$, while the critic estimates a value function $V_\phi(s)$ or the action-value function $Q_\phi(s, a)$, as illustrated in Fig. 2 (d). The critic provides feedback to the actor on how good the selected actions are, helping to stabilize the policy learning process. The objective of the actor is to maximize the expected return $J(\theta)$ defined as:

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right] \quad (16)$$

where γ is the discount factor and r_t represents the reward at time step t . The policy of the actor is updated by following the gradient of the expected return computed as:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) \cdot \delta_t] \quad (17)$$

where $\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$ is the temporal difference error. The critic updates the value function $V_\phi(s)$ by minimizing the TD error, thus improving its estimation over time and guiding the actor's policy updates.

Advantage Actor-Critic (A2C) is a synchronous variant of the standard AC architecture [61]. In A2C, the advantage function $A(s, a) = Q(s, a) - V(s)$ is explicitly calculated to reduce the bias introduced by the value estimation and to stabilize the updates. By focusing on the advantage rather than raw Q-values or returns, A2C improves the convergence speed and the stability of the learning process. The policy gradient in A2C is computed as:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) (Q(s, a) - V(s))] \quad (18)$$

A2C allows multiple workers to gather experience simultaneously to ensure that updates from each worker contribute to a stable and consistent policy. Asynchronous Advantage

Actor-Critic (A3C) builds on this idea but introduces an asynchronous framework where multiple workers interact with their environments independently. In A3C, each worker runs in its own instance of the environment and collects data independently and updates a shared global model. This setup leads to faster updates and reduces the likelihood of workers converging to suboptimal policies, which is a common issue in synchronous methods. By processing their local experiences and updating the global policy and value networks at different times, workers in A3C promote more robust learning and minimize the risk of overfitting to any single trajectory or environment. This results in faster learning speeds and more stable convergence.

2) *Proximal Policy Optimization (PPO)*: PPO presents a solution to the challenges inherent in policy gradient methods, particularly the risk of unstable or excessively large policy updates. Traditional approaches to policy optimization frequently experience abrupt policy changes after updates, making the learning process more challenging. PPO counteracts this by limiting how much the updated policy can diverge from the previous one using a trust region constraint. This is done via a surrogate objective function that controls the scale of policy changes. At its core, PPO relies on optimizing a clipped objective function, which ensures that policy updates remain within a specified range and do not destabilize the training process. The objective function in PPO is expressed as:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)] \quad (19)$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio between the new and old policies, A_t is the advantage function and ϵ is a small hyperparameter that controls the allowable change to the policy. The clipping function ensures that the policy ratio $r_t(\theta)$ stays within a bounded range $[1 - \epsilon, 1 + \epsilon]$, thus preventing large updates that could destabilize training.

PPO can be implemented in two different ways, either as PPO-Clip or PPO-Penalty. The most common form PPO-Clip applies the clipped objective as shown above, while PPO-Penalty uses a penalty term to constrain policy changes by minimizing the KL-divergence between the new and old policies:

$$L^{\text{PEN}}(\theta) = \mathbb{E}_t [L(\theta) - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot|s_t)||\pi_{\theta}(\cdot|s_t)]] \quad (20)$$

where β is a penalty coefficient, and the KL-divergence term $\text{KL}[\cdot||\cdot]$ ensures that the new policy does not deviate significantly from the old one.

3) *Deterministic Policy Gradient Methods (DPG)*: DPG methods are designed for continuous action spaces by directly learning a deterministic policy $\mu_{\theta}(s)$ which selects a specific action for each state. This approach avoids the variance introduced by sampling from a stochastic policy. The goal of DPG is to maximize the expected cumulative reward, formulated as:

$$J(\theta) = \mathbb{E}_{s_0 \sim p(s)} \left[\sum_{t=0}^T \gamma^t r(s_t, \mu_{\theta}(s_t)) \right] \quad (21)$$

DPG computes the policy gradient using the deterministic policy and the action-value function $Q(s, a)$, without requiring action sampling. The policy gradient is expressed as:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s \sim \rho^{\mu}} \left[\nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \Big|_{a=\mu_{\theta}(s)} \right] \quad (22)$$

DPG leverages experience replay and target networks to enhance learning stability. While effective in continuous domains, DPG can be prone to local optima due to insufficient exploration. Extensions like Deep Deterministic Policy Gradient (DDPG) address this limitation by incorporating noise-based exploration and lead to more robust performance.

In DDPG, the actor network learns a deterministic policy $\mu_{\theta}(s)$ while the critic network estimates the Q-value function $Q(s, a)$. To ensure exploration, DDPG adds noise to the policy during training typically through an Ornstein-Uhlenbeck process [62]. The critic network is trained using experience replay and target networks (see Equation 10), while the actor network is updated using the deterministic policy gradient (see Equation 22). These elements help stabilize learning and reduce the variance introduced by noisy gradients.

4) *Twin Delayed Deep Deterministic Policy Gradient (TD3)*: While DDPG performs well in continuous control tasks, it may still experience overestimated Q-values and result in suboptimal performance. Overestimation can lead to overly optimistic policy updates and degraded performance especially in complex continuous control tasks. TD3 introduces three main improvements to mitigate these issues: clipped double Q-learning, delayed policy updates, and target policy smoothing.

First, TD3 uses clipped double Q-learning to reduce overestimation. Instead of relying on a single Q-network, TD3 uses two Q-networks and updates the policy using the minimum Q-value between the two networks:

$$y = r + \gamma \min_{i=1,2} Q_{\theta_i}(s', \mu_{\theta'}(s')) \quad (23)$$

where $Q_{\theta_1}(s, a)$ and $Q_{\theta_2}(s, a)$ represent the estimates from the two critic networks.

The second improvement is delayed policy updates, which decouples the frequency of updates between the actor and critic networks. In TD3, the actor is updated less frequently than the critics, allowing the value function to stabilize before the policy is adjusted. This ensures that the policy is being updated based on more accurate value estimates, further reducing the risk of unstable updates.

Finally, target policy smoothing adds noise to the target actions used in the Q-value updates to prevent the policy from overfitting to narrow peaks in the value function. By applying small and clipped noise to the target actions, TD3 ensures that the policy remains robust to small errors in the value function approximation.

C. Multi-Agent DRL (MADRL) Methods

In cloud environments, agents are not isolated entities but interact with others. MARL extends traditional RL methods to such environments, where agents must learn policies in the presence of other agents whose behaviors also evolve during training [63] [64]. The dynamic nature adds complexity because agents must consider the strategies of others, which can range from cooperative to adversarial. Depending on the relationships between agents, multi-agent settings are typically classified as cooperative, competitive, or mixed, each requiring specific learning strategies and methods.

1) *Cooperative Setting*: In the cooperative setting, multiple agents share a common objective and their goal is to optimize a joint reward function, as depicted in Fig. 2 (e). This is typically formulated as a Decentralized Partially Observable Markov Decision Process (Dec-POMDP), where agents operate under partial observability and need to collaborate based on local observations to achieve a global goal [65].

A Dec-POMDP is defined by a tuple $\langle I, S, \{A_i\}_{i \in I}, P, R, \{O_i\}_{i \in I}, \gamma \rangle$. In this framework, I denotes the set of agents involved in the decision-making process, while S represents the set of global states in which the system can exist. Each agent $i \in I$ is associated with a distinct action set A_i , allowing it to select from a range of possible actions. The joint actions of all agents are denoted by $\mathbf{a} = (a_1, a_2, \dots, a_N)$, where each $a_i \in A_i$ corresponds to the action chosen by agent i . The state transition function denoted $P(s'|s, \mathbf{a})$, models the probability of transitioning from state s to a new state s' given the joint action \mathbf{a} . A shared reward function $R(s, \mathbf{a})$ determines the immediate reward based on the current state s and the joint actions of the agents. Each agent i receives observations from the environment, with the observation function O_i mapping the global state and joint actions to an observation specific to the agent. Finally, the discount factor γ balances future rewards against immediate ones, shaping the optimization of agent strategies.

The objective in cooperative MARL is to maximize the cumulative reward across all agents:

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \mathbf{a}_t) \right] \quad (24)$$

where $\pi = \{\pi_1, \pi_2, \dots, \pi_N\}$ represents the joint policy of all agents. Given the complexity of coordinating decentralized

agents, common techniques in this setting involve centralized training with decentralized execution (CTDE). This framework enables agents to learn jointly using shared information during training, but act independently during execution [66].

A prominent example of cooperative MARL is Multi-Agent Deep Deterministic Policy Gradient (MADDPG), which extends DDPG to multi-agent settings by introducing a shared critic that considers the joint actions and states of all agents [67]. With the actor decentralized, each agent is capable of learning its own policy autonomously:

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E} [\nabla_{\theta_i} \log \pi_{\theta_i}(a_i | o_i) Q_{\theta_i}(\mathbf{s}, \mathbf{a})] \quad (25)$$

where $\mathbf{a} = (a_1, a_2, \dots, a_N)$ and Q_{θ_i} is the centralized Q-function. The decentralized nature of the actor allows each agent to independently acquire its policy, which is particularly useful in scenarios such as multi-server or multi-cloud task parallelization.

2) *Competitive Setting*: In the competitive setting, agents have conflicting objectives, where each agent seeks to maximize its own reward often at the expense of others [68]. The corresponding architecture is presented in Fig. 2 (f). This is commonly modeled as a zero-sum game where the sum of the rewards for all agents at any time step is zero. The formal objective for agent i in a competitive setting is:

$$J_i(\pi_i) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \right] \quad (26)$$

subject to $\sum_{i=1}^N r_i(s, a) = 0$ for all s and a

Furthermore, in competitive environments characterized by adversarial agent actions, a central focus is the determination of the Nash equilibrium, a set of strategies from which no agent has an incentive to unilaterally deviate [69]. Competitive multi-agent DRL uses techniques such as policy gradient methods to optimize agents' policies in adversarial environments. A notable strategy is Self-Play where agents train by competing against themselves [70]. Specifically, the agents continuously update their policies by playing against increasingly skilled versions of themselves, leading to higher-level strategies over time. This approach has been successful in applications such as cloud resource allocation and scheduling, where agents learn complex strategies for efficient task distribution in highly dynamic environments. In general, adversarial training and min-max optimization are frequently employed in competitive MARL settings [71], where agents aim to maximize their own resource usage or task completion while minimizing the efficiency of competing agents:

$$\min_{\pi_1} \max_{\pi_2} \mathbb{E} \left[\sum_{t=0}^T \gamma^t R(s_t, \pi_1(s_t), \pi_2(s_t)) \right] \quad (27)$$

this competitive framework applies to multi-cloud task scheduling, load balancing, and resource allocation, where agents compete to optimize resource usage, dynamically adapting to workload or network fluctuations.

3) *Mixed Setting*: In the mixed setting, agents must navigate a balance between cooperation and competition, as their objectives may align with some agents while conflicting with others [72]. As a result, the reward function is neither purely shared nor purely adversarial, leading to the need for multi-objective optimization. For each agent i , the reward function may combine both individual and joint rewards:

$$r_i(s, \mathbf{a}) = \alpha_i r_{\text{shared}}(s, \mathbf{a}) + (1 - \alpha_i) r_i(s, a_i) \quad (28)$$

where α_i is a weight balancing the shared and individual rewards. Such scenarios frequently arise in real-world applications, where agents must balance both cooperation and competition depending on the specific task at hand. For example, in resource management, agents may collaborate to maximize overall system efficiency but compete for limited resources. The challenge in mixed settings is to dynamically adjust between cooperation and competition depending on the context, making it a highly challenging and significant area of study in MARL.

D. Advanced RL Methods

1) *Hybrid Heuristic and DRL*: The integration of heuristic methods with DRL offers a powerful framework for addressing complex optimization problems, leveraging the unique strengths of both approaches. Heuristic algorithms fundamentally search strategies driven by rule-based or experience-informed guidance, excel at swiftly generating approximate solutions within large or high-dimensional search spaces. These algorithms are highly efficient, flexible, and computationally economical, making them well-suited for providing rapid and near-optimal solutions. In contrast, DRL methods excel in learning optimal policies through sequential decision-making and are capable of continuous refinement based on cumulative feedback. The collaboration between heuristic algorithms and DRL can be established through several cooperative paradigms, each contributing to different stages of the problem-solving process.

Preliminary Solution Generation and Problem Decomposition: Before the DRL agent begins its training, heuristic algorithms can provide a initialization by decomposing the task or generating an approximate solution. In scenarios such as cloud workflow scheduling, heuristics can segment tasks based on resource requirements or workflow dependencies, effectively reducing the search space the DRL agent must navigate. This pre-processing phase enables DRL to initiate from a feasible solution space, accelerating the training process and reducing the exploration burden.

Enhanced Action Space Exploration: During the DRL training phase, heuristic methods can serve as auxiliary tools, guiding the agent's exploration within the action space. Traditional exploration strategies such as random sampling, often result in inefficient search trajectories particularly in high-dimensional environments or those with sparse rewards. Heuristic algorithms can mitigate these challenges by directing the agent toward regions of the action space with high potential rewards, thus prioritizing valuable paths and enhancing the efficiency of exploration. This heuristic-guided approach not

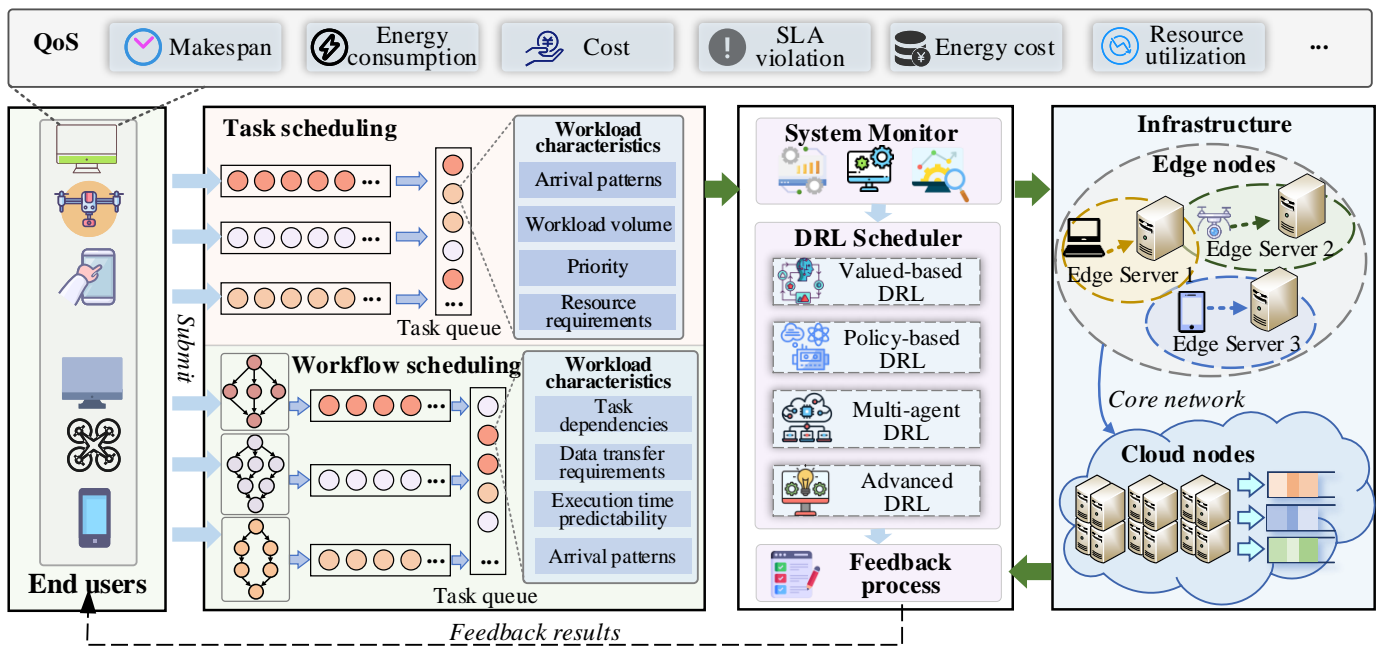


Fig. 3: The system architecture of job scheduling in cloud computing

only accelerates training but also increases the likelihood of discovering high-value solutions in complex, reward-scarce environments.

Targeted Adjustment and Solution Enhancement: Following the initial DRL optimization, heuristic algorithms can further refine the solutions generated. While DRL is effective at identifying broadly optimal solutions within complex spaces, it may lack the precision necessary for detailed improvements in specific regions of the solution space. Here, heuristic techniques can be applied to make targeted adjustments, leveraging their adaptability to enhance specific task arrangements or resource allocations. This staged optimization, in which DRL first generates an initial broad solution followed by heuristic-driven local refinements, leads to superior solution quality and greater adaptability to task-specific constraints.

In summary, the hybrid integration of heuristic methods and DRL establishes a layered, adaptable optimization approach, allowing agents to navigate and solve complex, multi-dimensional tasks more effectively. This synergy is especially advantageous in large-scale applications, where it accelerates convergence, enhances solution quality, and provides a versatile framework adaptable to diverse problem landscapes.

2) *Quantum Reinforcement Learning (QRL)*: With the rapid advancements in quantum technology, the advantages of applying QRL for optimization tasks are increasingly evident. One of the central techniques in QRL is the use of Quantum neural networks (QNNs), which can effectively embed DRL's state space into a quantum computational framework. Due to the characteristics of quantum circuits, a single state can often be embedded multiple times within the QNN, thus allowing for a richer representation and more diverse exploration. A key feature of QNNs in this context is the quantum variational approach, where quantum variational gates and entanglement gates replace the traditional neurons used in classical neural

networks, as shown in Fig. 2 (c). This shift significantly improves the efficiency of trainable parameters, as quantum gates allow for parallel computation and entanglement, enabling more complex state relationships to be represented with fewer resources.

Mathematically, the variational embedding can be expressed as a parameterized quantum circuit $U(\theta)$ acting on a quantum state $|\psi\rangle$ [73]. For example, an initial quantum state can be transformed by a sequence of variational gates:

$$|\psi(\theta)\rangle = U(\theta)|0\rangle^{\otimes n} \quad (29)$$

where θ represents the trainable parameters across the circuit, and $|0\rangle^{\otimes n}$ is the initial quantum state. Through repeated application of entanglement and rotation gates, such as CNOT and parameterized rotation gates $R(\theta)$, the QNN efficiently learns complex mappings of the DRL state space. Each layer in the QNN introduces entanglement across qubits, thereby capturing correlations within the data that are not easily accessible to classical neural networks.

Compared to traditional DNNs in DRL, QNNs in QRL can achieve comparable or even superior performance with significantly fewer parameters. Variational gates enable this efficiency by leveraging quantum entanglement and superposition, allowing QNNs to capture complex state relationships without the heavy computational load of DNNs. This efficiency makes QRL a promising approach as quantum technology advances [74].

IV. DRL-BASED CLOUD JOB SCHEDULING

In this section, we review existing works on DRL-based task scheduling and workflow scheduling, which form the core of job scheduling in cloud computing. Task scheduling involves assigning individual tasks to available resources

TABLE II: Summary of typical works on DRL-based task scheduling

Reference	Method	Others	Environment	Optimization objective				Others
				Makespan	Cost	Energy consumption		
[75]	DQN	-	Cloud computing	✓	-	✓	SLA violation	
[76]	DQN	-	Cloud computing	✓	-	-	Energy cost	
[11]	DQN	-	Cloud computing	✓	-	-	Energy cost	
[77]	DQN	Online learning	Cloud computing	✓	-	✓	Resource utilization	
[31]	DQN	-	Cloud computing	✓	✓	-	-	
[78]	DQN variants	-	Cloud computing	✓	-	✓	-	
[79]	DQN variants	Quantile regression	Cloud computing	✓	-	✓	-	
[80]	DDQN	Greedy optimization	Cloud computing	✓	✓	-	-	
[81]	DDQN	-	Cloud computing	✓	-	-	Task completion rate	
[82]	DQN	-	Edge computing	✓	-	-	Users served by base stations	
[83]	DQN	Online learning	Edge computing	✓	-	-	Avoiding severe task starvation	
[84]	DDQN	PSO	Edge computing	✓	-	-	-	
[85]	DDQN	-	Edge computing	-	-	✓	Training efficiency	
[86]	PPO	-	Cloud computing	✓	-	✓	-	
[87]	PPO	-	Cloud computing	-	-	✓	-	
[88]	PPO	-	Cloud computing	✓	✓	-	Renewable energy utilization; deadline violation	
[89]	DDPG	-	Cloud computing	✓	-	-	CPU utilization standard deviation	
[90]	DDPG	-	Cloud computing	✓	✓	-	-	
[91]	SAC	GCN	Edge-Cloud collaboration	✓	-	✓	-	
[6]	A2C	-	Edge-Cloud collaboration	-	-	-	Server resource utilization; task rejection rate	
[92]	A3C	-	Edge-Cloud collaboration	✓	✓	✓	SLA Violations	
[93]	DDPG	-	Edge computing	-	✓	-	System utility	
[94]	DDPG	-	Edge-Cloud collaboration	-	-	-	Load forward to the cloud server	
[95]	MADRL	-	Cloud computing	-	-	-	Reliability; network communication overhead; resource utilization	
[96]	MADRL	-	Cloud computing	-	✓	✓	Fairness	
[97]	MADRL	-	Cloud computing	✓	-	-	Resource utilization	
[98]	MADRL	-	Edge computing	✓	-	✓	Throughput rate	
[99]	D3RQN	-	Edge cloud	✓	-	✓	Utility	
[100]	MADRL	-	Edge computing	-	✓	-	QoS satisfaction rates	
[101]	MADRL	SA	Edge-Cloud collaboration	✓	-	✓	Throughput rate	
[102]	MADRL	-	Edge-Cloud collaboration	✓	-	✓	-	
[103]	MAPPO	GNN	Edge computing	✓	-	-	Resource efficiency	
[104]	Multi-Agent PPO	Meta learning	Edge computing	-	-	-	Resource utilization	
[105]	MADDPG	-	Edge computing	✓	-	✓	-	
[106]	D3QN	Adversarial imitation learning	Cloud computing	✓	-	✓	-	
[107]	DQN	LSTM	Cloud computing	-	-	-	CPU usage cost; RAM memory usage cost	
[108]	Distributional RL	Quantile regression	Cloud computing	✓	-	-	Load balancing; success rate	
[109]	Double-level DRL	Interactive training strategy	Cloud computing	-	✓	-	Computation efficiency	
[110]	DDQN	Representation learning	Edge computing	-	-	✓	SLA	
[111]	A3C	Multi-task learning	Edge-Cloud collaboration	✓	✓	-	Load imbalance value	

to optimize system performance, while workflow scheduling manages interdependent tasks that must execute in a specific order, focusing on task coordination and dependency resolution. As illustrated in Fig. 3, the scheduling process starts with user requests submitted to the cloud platform. DRL-based approaches allow the scheduler to dynamically analyze workload characteristics and resource availability, enabling efficient task allocation while maintaining adherence to QoS requirements.

A. DRL-Based Task Scheduling

Through continuous interaction with the environment and learning from experience, DRL has demonstrated exceptional effectiveness in addressing task scheduling challenges in cloud and edge computing environments [75], [76]. Below, we review key works that employ DRL for task scheduling, with a summarized overview presented in Table II.

1) *Valued-Based DRL for Task Scheduling*: Value-based DRL techniques, such as DQN and its variants, have demonstrated robust decision-making capabilities by estimating state-action value functions for task scheduling. Specifically, these methods focus on learning an optimal policy by approximating the action-value function, which captures the expected reward of selecting specific task scheduling actions, such as assigning tasks to resources or determining execution orders, across various scheduling scenarios [20]. Extensive studies have investigated the application of valued-based DRL to

enhance task scheduling strategies in cloud computing, achieving substantial improvements in task execution efficiency, and resource utilization. As a foundational approach within valued-based DRL, DQN has been effectively utilized to address the inherent challenges of task scheduling including dynamic workload characteristics [11], workload balancing [77], and the unpredictable nature of cloud environments [31]. For instance, the work [75] successfully applies DQN to minimize energy consumption, makespan and SLA. Likewise, the work [76] introduces a DQN-based method that employs a two-stage processor, where the first stage assigns tasks to appropriate server clusters and the second stage selects specific servers within those clusters for execution. In the work [11], a DQN-driven energy-aware scheduling approach is introduced, which dynamically allocates tasks to optimal VMs based on real-time workload characteristics and resource availability, thereby reducing energy consumption while ensuring QoS. To further adapt to dynamic workloads, the work [77] develops an adaptive DQN framework, which modifies the discount factor γ to optimize energy usage in response to workload fluctuations. Furthermore, the work [31] presents a DQN-based preemptive method aimed at optimizing job execution cost and response time.

While DQN has demonstrated effectiveness in task scheduling, various DQN variants have also been extensively employed to further improve scheduling quality in cloud environments. These variants have been designed to address specific

limitations of conventional DQN approaches, such as Q-value overestimation and the need for adaptation to specialized tasks. For example, the work [78] presents a parameterized action space-based DQN (PADQN) framework for jointly optimizing task dispatching and cooling regulation. PADQN specifically addresses the hybrid action space problem by concurrently managing the discrete actions involved in task dispatching and the continuous actions required for cooling regulation. In another work [79], a quantile regression DQN (QR-DQN) network is employed to determine an optimal long-term scheduling strategy, effectively addressing the uncertainties in cloud workload patterns. Given that conventional DQN approaches often suffer from Q-values overestimation, DDQN has been adopted to mitigate this issue. For instance, the work [80] combines DDQN with the greedy optimization strategy for online task scheduling, where DDQN handles task assignments and the greedy algorithm focuses on task execution. Similarly, the work [81] applies DDQN to improve task completion rates while simultaneously reducing average response time, demonstrating its effectiveness in handling complex scheduling scenarios under dynamic cloud conditions.

Value-based DRL methods have also been effectively applied in edge computing environments to address challenges such as resource constraints, fluctuating network conditions, energy consumption, and communication delays. For instance, the work [82] employs DQN to efficiently schedule tasks to mobile edge computing (MEC) servers, optimizing both delay and bandwidth utilization. Similarly, the work [83] presents a DQN-based method designed to mitigate task delays and alleviate task starvation caused by fluctuating network conditions and unbalanced server loads in edge computing environments. To address the issue of Q-value overestimation in traditional DQN, DDQN is introduced to improve scheduling accuracy in these environments. For example, the work [84] integrates PSO with DDQN to minimize task scheduling time in edge computing, thereby improving resource utilization efficiency and overall QoS. Additionally, the work [85] leverages DDQN to achieve energy-efficient task scheduling at the network edge by incorporating dynamic voltage and frequency scaling (DVFS) mechanisms, where the evaluation network calculates Q-values for various DVFS configurations, and the target network produces estimated Q-values to guide parameter optimization.

2) *Policy-Based DRL for Task Scheduling*: Value-based DRL methods have achieved notable progress in task scheduling. However, their dependence on value function estimation often restricts their performance in complex and dynamic scenarios, resulting in challenges such as value overestimation and suboptimal convergence. To address these limitations, policy-based DRL methods offer a promising alternative by directly optimizing scheduling policies. Among these, PPO has shown considerable potential in tackling task scheduling challenges within cloud environments. For example, in the work [86], an enhanced PPO algorithm is proposed, integrating priority rules to manage task scheduling with deadlines, particularly in scenarios characterized by random task arrivals and renewable energy integration. This approach aims to minimize both service costs and latency, demonstrating the adaptability of

PPO to dynamic and unpredictable scheduling environments. Moreover, the work [87] utilizes PPO to train agents, which enhances energy efficiency and system performance for cloud task scheduling. In hybrid environment, PPO has proven to be a powerful tool for optimizing task scheduling. The work [88] introduces an advanced task scheduling framework based on PPO, designed to maximize renewable energy utilization while strictly adhering to deadline requirements in hybrid cloud environments. This framework leverages PPO to enable real-time workload shifting and decision-making, optimizing key metrics such as operational cost, makespan, and resource utilization. In addition to PPO, several studies have explored DDPG-based techniques for task scheduling in cloud computing environment. For instance, in the work [89], a DDPG-based task scheduling technique is introduced to tackle load balancing and SLA assurance, with a particular focus on optimizing energy consumption in data centers. Expanding on this, the work in [90] extends DDPG to handle large-scale and heterogeneous cloud workloads, aiming to optimize both response time and cost. This approach leverages a dual reward model to enable agents to effectively learn optimal scheduling policies under complex workload conditions.

Various policy-based DRL approaches have been proposed to address complexities for task scheduling in edge computing effectively. For example, the work [91] leverages the Soft AC algorithm to optimize task scheduling within edge-cloud environments. Additionally, this approach integrates Graph Convolutional Networks (GCN) to capture complex dependencies among tasks through graph-based representations, which are essential for modeling interdependencies in edge-cloud environments. Moreover, the work [6] employs the A2C algorithm to dynamically schedule tasks in response to fluctuations in edge environments, with the objective of maximizing server resource utilization while minimizing task rejection rates. In another notable approach, the work [92] introduces a task scheduling framework based on the A3C algorithm. This framework incorporates Residual Recurrent Neural Networks (R2N2) to update model parameters in real time, enabling effective adaptation to the stochastic nature of edge-cloud environments. Several works have explored the potential of DDPG in addressing task scheduling challenges in edge computing environments. For example, the work [93] utilizes a DDPG-based algorithm that adjusts network structures to accommodate the discrete action space of edge environments. This approach is specifically designed to enhance scheduling efficiency by maximizing system utility while minimizing cumulative operational costs, effectively meeting the demands of dynamic and resource-constrained edge computing settings. Furthermore, the work [94] presents an advanced dynamic task scheduling framework based on the DDPG algorithm, tailored for edge-cloud Internet of Things (IoT) systems. This framework proficiently manages service migration while meeting stringent latency constraints, ensuring optimal performance in highly dynamic settings.

3) *Multi-Agent DRL for Task Scheduling*: The application of MADRL in task scheduling has gained significant attention in recent years [95], [99]. MADRL leverages the collaborative capabilities of multiple agents, facilitating dynamic task

scheduling while overcoming the scalability and complexity issues that often hinder traditional single-agent approaches [97]. In cloud environment, for instance, the work [95] utilizes a MADRL framework for task scheduling, focusing on mitigating total system failure risk and resolving the single point of failure challenge. Additionally, the work [96] presents a cloud-assisted MADRL scheduling framework to optimize task scheduling and energy management in a Unmanned Aerial Vehicle (UAV) charging network. This framework employs a centralized orchestration manager to coordinate energy sharing and scheduling, achieving efficient and adaptive energy distribution across multiple charging stations. Moreover, the work [97] proposes a general-purpose MADRL framework designed to learn optimal collaborative task scheduling policies. This framework demonstrates the ability to adapt to varying workload demands and effectively manage resource allocation.

Beyond cloud environments, MADRL has been increasingly applied to edge computing, where challenges such as resource heterogeneity, dynamic network conditions, and fluctuating task requirements are prevalent. In these environments, MADRL enables distributed decision-making among edge nodes, improving system efficiency and responsiveness. For instance, the work [98] introduces a Value-Decomposition Multi-Agent DQN (VD-MADQN) for real-time scheduling of user requests within edge networks. The proposed approach enables edge nodes to independently make scheduling decisions based on localized information while leveraging centralized training to foster cooperative strategies. Additionally, the work [99] formulates task scheduling in serverless edge computing network as a partially observable stochastic game (POSG). By employing a dueling double deep recurrent Q-network (D3RQN) algorithm, this MADRL framework empowers each edge computing node to autonomously make scheduling decisions, effectively utilizing local observations while aligning with global optimization objectives.

In addition to value-based approaches, recent advancements in MADRL have incorporated policy-based methods to further optimize task scheduling in edge computing. For example, the work [100] proposes a MADRL-based framework built on AC, tailored for distributed transmission within collaborative cloud-edge environments. This approach focuses on joint user scheduling and beam selection, aiming to minimize long-term network delay while maintaining adherence to QoS constraints. In large-scale industrial IoT applications, the work [101] introduces a collaborative MADRL framework using the A3C algorithm. This framework allows agents to adapt dynamically to changing conditions, fostering cooperation among agents, and improving both convergence and system stability. Furthermore, the work [102] develops a competitive multi-agent Attention-Communication AC (MA3C) to support diverse task types within cloud-edge-end collaborative systems. This approach utilizes attention mechanisms for agents to focus on relevant information from other agents in a partially observable environment, thereby enhancing load balancing and overall system efficiency. Similarly, the work [103] employs a multi-agent PPO (MAPPO) framework for task scheduling in distributed edge computing networks. This framework integrates specific adaptations for edge

environments, including state abstraction via heterogeneous graph attention networks (HAN) to capture complex inter-agent semantics and action decomposition for efficient task selection. To address the nonstationarity challenges inherent in heterogeneous edge computing, the work [104] develops a multi-agent meta-PPO algorithm. By leveraging meta-learning techniques, this approach accelerates convergence and improves overall system efficiency and stability under dynamic conditions. Lastly, the work [105] provides a comprehensive analysis of two scenarios. In single-edge settings, a DRL-based framework is employed for collaborative task scheduling, while in multi-edge scenarios, the MADDPG algorithm is applied to minimize energy consumption and latency, ensuring efficient resource utilization across edge nodes.

4) *Advanced DRL for Task Scheduling*: To overcome the inherent complexities for task scheduling in cloud computing environments, recent research has incorporated advanced techniques into DRL frameworks, pushing the boundaries of task scheduling optimization. For instance, the work [106] proposes a hybrid framework that integrates adversarial imitation learning with the Dueling Double Deep Q-Network (D3QN) algorithm for cloud task scheduling. By utilizing adversarial imitation learning to store high-reward job trajectories as expert demonstrations, this approach directly guides the DRL agent, enhancing both policy optimization and scheduling performance. The work [107] presents a DRL framework integrated with LSTM networks to optimize VM scheduling in big data analytics. The LSTM network captures long-term dependencies between tasks and resource demands, enhancing scheduling efficiency and reducing execution costs. The work [108] proposes a distributional RL-based approach for load balancing in cloud computing environments, with a focus on batch task scheduling. By employing quantile regression, this framework effectively distributes computational loads across resources, adapting dynamically to fluctuations in job requirements and cluster states. Additionally, the work [109] introduces a Double-Level DRL approach that includes an Interactive Training Strategy (ITS), designed to boost adaptability and scalability.

As edge computing environments expand in scale and complexity, advanced DRL techniques are increasingly applied to address the demands of dynamic, high-dimensional scheduling tasks. For instance, the work [110] proposes a task scheduling framework that combines DDQN with representation learning to handle the complexities of dynamic edge environments. By reducing the dimensionality of nodes and tasks, the representation learning component enables efficient high-dimensional data processing, thereby enhancing the speed and accuracy of DRL-based decision-making in edge computing environments. Additionally, the work [84] improves DDQN efficiency by incorporating a pre-training phase with PSO, which provides a near-optimal initialization for task scheduling. This hybrid approach accelerates the initial learning phase of DDQN and addresses overestimation issues in DQNs by decoupling action selection from target Q-value computation. In a parallel direction, the work [111] introduces a scalable multi-task DRL framework for parallel task scheduling, which leverages multi-task learning to efficiently manage high-dimensional

action spaces and concurrently optimize multiple tasks. This approach reduces computational overhead while enhancing resource utilization in edge environments.

B. DRL-Based Workflow Scheduling

Workflow applications are increasingly prevalent in cloud environments due to their effectiveness in addressing complex challenges in scientific research, such as astronomy, bioinformatics, and seismology [154]. A workflow application typically consists of multiple computational tasks organized as a DAG, where nodes represent individual tasks and edges indicate dependencies among these tasks. The inherent complexity of such a structure makes workflow scheduling an NP-hard combinatorial optimization problem. Recent research has reformulated workflow scheduling as a discrete-time control problem, leveraging DRL techniques to develop more adaptive and scalable scheduling frameworks [112], [113], [114]. Below, we provide a detailed analysis of these studies and summarize them in Table III.

1) *Valued-Based DRL for Workflow Scheduling*: Valued-based DRL methods, represented by DQN, have demonstrated their effectiveness in addressing workflow scheduling challenges within cloud computing environments. For example, the work [112] employs DQN to handle DAG-based tasks in a cloud computing environment, with a primary focus on reducing makespan variance and improving load balance. Furthermore, DQN has been successfully applied in real-time workflow scheduling. In a dynamic environment with multiple real-time workflows, the work [113] applies DQN to assign each task in workflows to a suitable VM, thereby minimizing workflow makespan and maximizing resource utilization. Similarly, the work [114] leverages DQN to optimize real-time workflow scheduling in cloud environments, aiming to reduce execution time and dynamically manage resource allocation. To further enhance the performance of DRL in workflow scheduling, heuristic algorithms have been integrated with DQN, resulting in hybrid models that accelerate convergence and mitigate the risk of local optima. For instance, the work [115] introduces a hybrid model combining DQN with SA for cost-sensitive workflow scheduling in cloud environments. In this approach, SA optimizes the task execution sequence, while DQN learns optimal scheduling policies in dynamic environments, ultimately improving resource utilization and reducing scheduling costs. Likewise, the work [116] integrates GA with DQN, where GA is responsible for optimizing VM execution plans through a global search strategy, effectively reducing the exploration space of DQN. This integration minimizes execution costs and response times, thereby enhancing scheduling efficiency.

Building on the success of DQN in cloud-based workflow scheduling, its variants, such as DDQN [117] and D3QN [119] have significantly expanded its decision-making capabilities, enabling more effective solutions to complex scheduling challenges. For instance, the work [117] proposes a DDQN-based scheduling framework that mitigates Q-value overestimation by decoupling action selection from evaluation, thus minimizing task completion time and resource consumption. Similarly,

the work [118] introduces a weighted DDQN approach that incorporates adaptive dynamic coefficients to balance Q-value overestimation in DQN with the underestimation in DDQN. This approach enables simultaneous optimization of execution time and cost. Moreover, the work [119] introduces a D3QN-based collaborative scheduling strategy for heterogeneous workflows in cloud. By incorporating dueling and double Q-learning mechanisms, this strategy effectively addresses Q-value overestimation while optimizing workflow makespan, cost, fairness, and continuity. Additionally, the work [120] incorporates a bias correction mechanism to further alleviate Q-value overestimation, resulting in optimized task execution times and balanced load distribution across multi-cloud environments.

Beyond cloud computing, recent advancement in valued-based DRL methods have driven the development of more sophisticated workflow scheduling frameworks tailored to the unique requirements of edge computing environments. For example, the work [121] presents a DQN-based workflow scheduling framework for dynamic edge-cloud workloads, utilizing Implicit Quantile Networks (IQN) to enhance robustness and Ape-X for distributed experience replay. Trained offline using a cluster simulator, this framework adaptively selects the most suitable scheduling strategies to improve resource allocation and minimize response times under fluctuating conditions. Moreover, the work [122] develops a DQN-based scheduling algorithm that aims to balance workload distribution, reduce service time, and minimize task failure rates by optimizing task allocation across edge and cloud resources. Within IoT applications, the work [123] utilizes DQN to develop a workflow-aided IoT paradigm that integrates intelligence edge computing. This method automates the handling of task dependencies, improving VM allocation and minimize network delay. To address the limitations of DQN, some studies have explored advanced DQN variants to enhance scheduling efficiency and robustness under dynamic edge-cloud conditions. For example, the work [124] presents a hierarchical workflow scheduling framework tailored for collaborative cloud-edge-end computing environments with a DDQN approach. Furthermore, the work [125] employs a D3QN-based method that addresses Q-value overestimation by leveraging two Q-networks updated in an alternating manner. This method effectively optimizes key performance metrics such as execution time, energy consumption, and operational costs within edge computing environments.

2) *Policy-Based DRL for Workflow Scheduling*: Policy-based DRL methods offer significant advantages for workflow scheduling by directly optimizing scheduling policies. These methods effectively manage dynamic resource variations and adapt to changes in task arrivals, improving workflow execution efficiency and resource utilization through timely policy updates. The AC framework in policy-based DRL has been effectively applied to tackle the challenges of workflow scheduling in cloud computing. For example, the work [126] presents an AC framework to manage the complexities of cloud workflow scheduling. This method integrates a P-Network model for task prioritization prediction, supplemented by a heuristic algorithm that facilitates server selection based on

TABLE III: Summary of typical works on DRL-based workflow scheduling

Reference	Method	Others	Environment	Optimization objective			
				Makespan	Cost	Energy consumption	Others
[112]	Deep Q-learning	-	Cloud computing	✓	-	-	Load balance
[113]	DQN	-	Cloud computing	✓	-	-	Resource utilization
[114]	DQN	SA	Cloud computing	✓	✓	-	Task failures; communication costs
[115]	DQN	SA	Cloud computing	✓	✓	-	-
[116]	DQN	GA	Cloud computing	✓	✓	-	-
[117]	DDQN	-	Cloud computing	✓	-	-	Resource usage
[118]	DDQN	-	Cloud computing	✓	✓	-	-
[119]	D3QN	-	Cloud computing	✓	✓	-	Fairness and continuity
[120]	DQN	-	Cloud computing	✓	-	-	Load balance
[121]	DQN	-	Edge-Cloud collaboration	✓	-	-	-
[122]	DQN	-	Edge-Cloud collaboration	✓	-	-	VM utilization; failed tasks rate
[123]	DQN	-	Edge-Cloud collaboration	✓	-	-	Network throughput
[124]	DDQN	LSTM	Edge-Cloud collaboration	✓	-	✓	Utility
[125]	D3QN	-	Edge computing	✓	✓	✓	-
[126]	AC	-	Cloud computing	✓	-	-	-
[127]	AC	LSTM	Cloud computing	✓	-	-	-
[128]	AC	-	Cloud computing	-	-	-	Bounded slowdown; resource utilization
[129]	A2C	-	Cloud computing	-	-	✓	Carbon emission
[130]	PPO	GCN	Cloud computing	✓	-	-	Job latency
[131]	PPO	GCN	Cloud computing	✓	-	-	Resource utilization
[132]	PPO	Self-attention mechanism	Cloud computing	-	✓	-	-
[133]	DDPG	-	Cloud computing	✓	-	-	-
[134]	AC	-	Edge-Cloud collaboration	✓	-	✓	-
[135]	AC	GCN	Edge-Cloud collaboration	✓	-	-	Energy cost; network traffic; load balance
[136]	A3C	-	Edge-Cloud collaboration	-	-	-	Task rejection rate; resource utilization
[137]	PPO	GNN	Edge-Cloud collaboration	✓	-	-	Migration time; energy cost
[138]	PPO	-	Edge-Cloud collaboration	✓	✓	-	Load balance
[139]	PPO	-	Edge computing	✓	-	-	Resource utilization
[140]	DDPG	-	Edge computing	✓	-	-	-
[141]	MADRL	-	Cloud computing	-	✓	✓	Load balance; resource utilization
[142]	MADRL	-	Cloud computing	✓	✓	-	-
[143]	MADDPG	-	Cloud computing	✓	-	✓	-
[144]	MADRL	-	Cloud computing	✓	✓	-	Execution interruptions
[145]	MADRL	-	Edge computing	✓	-	✓	Success rate; resource utilization
[146]	MADRL	-	Edge computing	✓	-	✓	Success rate; load balance
[147]	MADRL	-	Edge computing	✓	-	-	-
[148]	Depth-First-Search Coalition RL	-	Cloud computing	✓	-	-	-
[149]	DQN	Transformer	Cloud computing	✓	✓	✓	-
[150]	DDQN	GNN	Cloud computing	✓	-	-	-
[151]	DQN	Primary backup	Edge computing	✓	-	-	-
[152]	DQN	QML	Edge-Cloud collaboration	✓	✓	-	Power utilization; violation of delay
[153]	DQN	Transformer	Edge-Cloud collaboration	✓	✓	✓	Weighted cost

task precedence and computational complexity considerations. Similarly, the work [127] proposes a dynamic workflows scheduling approach based on the AC model, with the objective of minimizing makespan. This approach utilizes an extended pointer network and the LSTM-based encoder to handle the dependency structure of workflows, enabling adaptive task allocation that improves efficiency under dynamic resource conditions. Additionally, the work [128] introduces the scheduling of DAG-based workflows using the AC method. This approach considers both the dependency structure within DAGs and the resource capacities of cloud data centers, aiming to improve scheduling decisions and optimize resource usage. Moreover, the work [129] proposes an eco-friendly A2C-based framework for workflow scheduling in federated cloud environment. The objective is to minimize energy consumption and carbon emissions by intelligently selecting data centers for workflow task execution.

The PPO algorithm has also been widely adopted to enhance workflow scheduling efficiency, advancing beyond the standard AC framework by introducing a stable policy optimization mechanism that constrains update magnitudes in cloud environment. For example, the work [130] proposes a PPO-based workflow scheduler for cluster, aiming to minimize job latency and makespan while prioritizing critical tasks. By leveraging the GCN to extract workflow features, the proposed scheduler optimizes resource allocation across multiple queues. Similarly, by incorporating Monte Carlo Tree

Search (MCTS) and GCN into PPO, the work [131] focuses on reducing workflow makespan while increasing resource utilization in the cloud environments. In heterogeneous cloud environments, the work [132] utilizes PPO for scheduling workflow, with the objective of minimizing instance costs and ensuring timely completion of all tasks. This approach incorporates a self-attention mechanism for effective feature extraction and employs a mask layer to prevent illegal actions, thereby enhancing the stability and convergence speed of the learning process. In addition to PPO, DDPG has also been explored for workflow scheduling. For example, the work [133] proposes a workflow scheduling policy based on the DDPG algorithm to minimize computation latency in distributed cloud computing systems.

The heterogeneity and dynamic characteristics of edge-cloud environments present additional challenges for workflow scheduling, further driving the adoption of policy-based DRL methods. For example, the work [134] proposes a hybrid AC model for optimizing workflow scheduling in edge-cloud environments, aiming to minimize energy consumption while ensuring timely task execution. This approach integrates multiple actor networks with a single critic network, effectively supporting hierarchical action spaces that distinguish between edge and cloud nodes. Expanding on similar principles, the work [135] proposes a framework based on the AC model to optimize the scheduling of componentized tasks within cloud-edge environments. This framework integrates a Graph Neural

Network (GNN)-enhanced DRL, utilizing a combination of GCN and Directed GCN to effectively capture and represent task and network graph information, optimizing multiple objectives such as system latency, energy cost, network traffic, and load balancing. Additionally, the work [136] utilizes the A3C algorithm for data-intensive workflow scheduling in a volunteer edge-cloud environment, considering workflow QoS requirements, security specifications, and the resource preferences of volunteer nodes (VNs). The A3C-based model employs both actor and critic networks within a parallel learning architecture involving multiple worker agents, enabling efficient adaptation to dynamic and heterogeneous environments.

PPO has also been effectively utilized to address workflow scheduling challenges in edge-cloud environments, offering robust solutions for managing dynamic and time-sensitive tasks. For example, the work [137] explores a PPO-based workflow scheduling in the edge-cloud computing environment with continuous task arrivals, leveraging a GNN-based workflow embedding to capture latent task dependency information. The framework introduces an intrinsic reward mechanism to provide immediate feedback and improve scheduling decisions dynamically, ultimately aiming to minimize makespan and enhance energy utilization. In addition, the work [138] proposes a PPO-based IoT application scheduling algorithm, which aims to optimize system load balancing and response time in edge and fog computing environments. To further enhance scheduling in time-sensitive and resource-constrained edge environments, the work [139] develops a workflow scheduling framework based on PPO, incorporating a self-critic mechanism for improved decision-making efficiency and transformer-based neural networks for better sub-task feature processing. Furthermore, the DDPG algorithm has been employed by the work [140] for workflow scheduling in the edge network, focusing on minimizing makespan. This framework integrates critical path analysis-based dynamic task sorting and a path quality indicator for multi-path routing, effectively coordinating computing resources and managing interdependent tasks.

3) *Multi-Agent DRL for Workflow Scheduling*: MADRL has demonstrated significant potential for tackling workflow scheduling challenges. In cloud environments, MADRL has been widely adopted for workflow scheduling, with the goal of solving complex scheduling challenges and improving overall system performance. For example, the work [141] proposes a workflow scheduling framework that leverages MADRL to manage multiple online scientific workflows. This framework integrates cooperative Q-learning with Markov game theory to optimize both resource provisioning and task scheduling, effectively balancing the objectives of reducing energy consumption, minimizing user costs, and distributing workloads evenly across resources for both users and cloud providers. Additionally, the work [142] proposes a MADRL framework based on DQN for multi-objective workflow scheduling in cloud environment. This framework focuses on optimizing workflow makespan and cost, dynamically adapting to different workflow types and VM configurations to achieve effective scheduling. Moreover, the work [143] presents a MADRL ap-

proach based on AC, specifically designed to optimize renewable energy usage in workflow scheduling across distributed cloud data centers. By employing a hierarchical approach, the framework features a global RL agent to assign tasks to data centers, while local RL agents assign tasks to individual nodes. This setup effectively manages the complexities of partial observability and distributed environments. Additionally, the work [144] utilizes a MADRL framework based on PPO to minimize workflow execution costs by optimizing the use of both preemptible and on-demand instances. This approach features multiple actor networks coordinated under the guidance of a centralized critic network, leveraging a hierarchical action space to ensure an optimal balance between resource utilization and operational efficiency in cloud environments.

In edge-cloud environments, the use of MADRL has demonstrated considerable potential in optimizing workflow scheduling under dynamic conditions. For instance, the work [145] presents a telemetry-aided cooperative MADRL framework based on DQN for workflow scheduling in edge computing environment. By utilizing telemetry data for real-time decision-making, this framework aims to improve resource allocation efficiency and workflow scheduling across edge servers and programmable switches, ultimately enhancing the performance of workflow execution under dynamic network conditions. Moreover, the work [146] integrates Digital Twin (DT) technology with Evolutionary Selection Multi-Agent Reinforcement Learning (ES-MARL) to address workflow scheduling challenges in large-scale MEC environments. In this approach, DT assists in maintaining up-to-date network states for centralized training, while distributed agents make local scheduling decisions. The use of the Multi-Agent Transformer (MAT) algorithm, combined with evolutionary selection, improves training efficiency and fosters effective collaboration among agents. In addition, the work [147] proposes a workflow scheduling approach based on MADRL for vehicular edge computing (VEC) environments. This method models the task offloading problem as a potential game and leverages DRL to enhance collaboration between vehicles and Roadside Units (RSUs), resulting in improved offloading decisions and resource utilization.

4) *Advanced DRL for Workflow Scheduling*: Recent advances in DRL have driven significant progress in workflow scheduling by enhancing workflow representations [150], integrating hybrid methods like quantum machine learning (QML) [152], and improving scheduling efficiency through sophisticated neural architectures such as Transformers [149], [153], etc. These advancements have transformed both cloud and edge computing environments. For cloud environment, advanced DRL approaches have introduced effective solutions for handling complex scheduling demands. For example, the work [148] presents an adaptive multi-workflow scheduling framework for cloud computing environments, employing a Depth-First-Search Coalition Reinforcement Learning (DF-SCRL) policy. This policy integrates physical machines (PMs) coalition formation with Q-learning to determine the optimal bundle of VM instances, leading to improved resource efficiency, reduced costs, and enhanced workflow reliability. Moreover, the work [149] extends the DRL by proposing

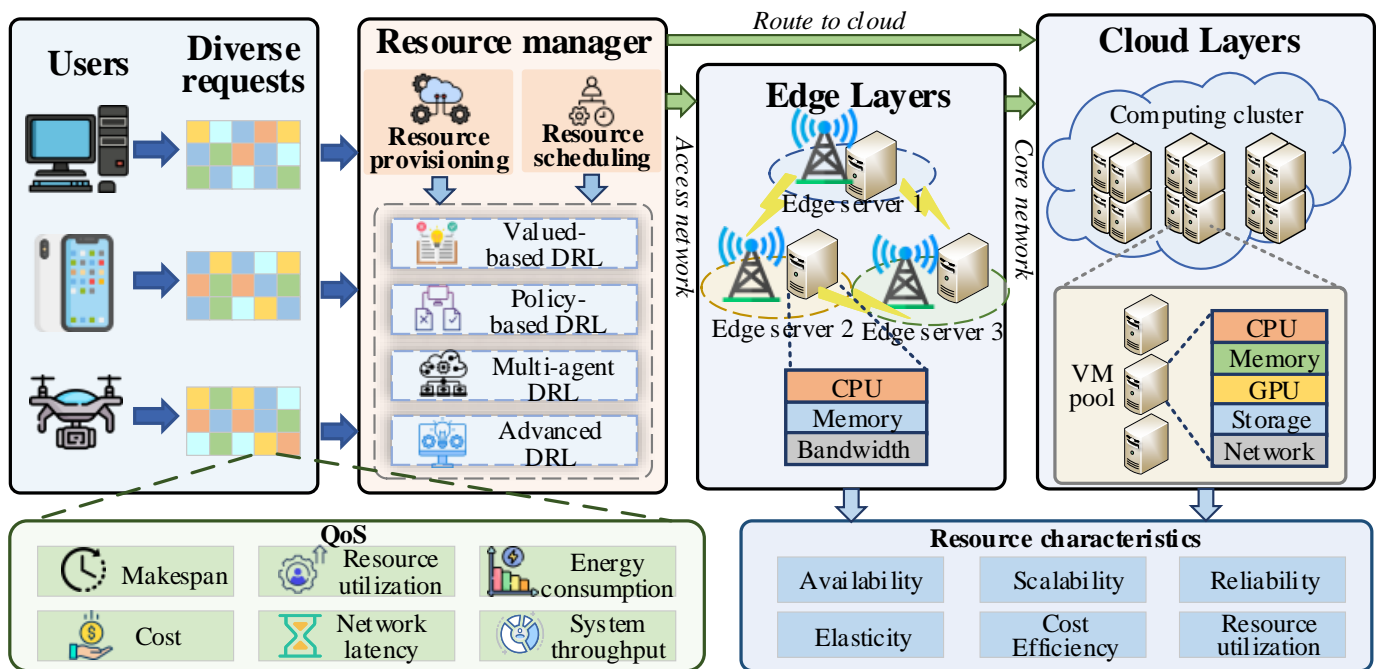


Fig. 4: The system architecture of resource management in cloud computing

a transformed-enhanced DQN approach for efficient large-scale dynamic workflow scheduling in heterogeneous cloud environments. By integrating transformer models, this approach effectively handles complex dependencies, uncertainties in task execution times, and dynamic resource availability, thereby optimizing system performance across diverse workflow scenarios. Additionally, the work [150] proposes a GNN-augmented DDQN framework for workflow scheduling in dynamic cloud environments. This framework employs a multi-head Graph Attention Network (GAT) to extract features of subtasks, accounting for both predecessor and successor relationships.

For edge computing, advanced DRL methods have addressed the unique challenges posed by distributed and resource-constrained environments. For example, the work [151] presents a fault-tolerant workflow scheduling method leveraging the Primary-Backup (PB) strategy in combination with DQN. This method enhances workflow reliability by mitigating resource and link failures, which are common in distributed edge environments. Moreover, the work [152] introduces a novel hybrid approach that integrates DQN with QML for workflow scheduling across edge, fog, and cloud layers. By utilizing QML principles such as superposition and entanglement, the proposed method effectively addresses challenges like power consumption, delay minimization, and service cost reduction in collaborative edge-cloud environments. Additionally, the work [153] introduces a distributed DRL framework enhanced by transformer for workflow scheduling in edge-cloud environments. By incorporating Prioritized Experience Replay (PER) and transformer layers into the DRL architecture, the framework effectively reduces high exploration costs and captures long-term dependencies among tasks, leading to better system performance.

V. DRL-BASED CLOUD RESOURCE MANAGEMENT

In this section, we review existing works on DRL-based resource provisioning and scheduling for cloud resource management. Resource provisioning focuses on allocating virtualized resources to meet user demands or workload requirements, while resource scheduling involves efficiently assigning these resources to tasks. DRL-based resource management has proven to be an adaptive and efficient solution for addressing the complexities of both provisioning and scheduling. By leveraging DRL techniques, resource managers dynamically optimize utilization, minimize operational overhead, and ensure compliance with QoS requirements across cloud and edge environments. Fig. 4 illustrates the architecture of a DRL-based resource management system, showcasing its flexibility in adapting to diverse and dynamic user demands.

A. DRL-Based Resource Provisioning

Resource provisioning is a critical aspect of resource management within cloud computing, with an emphasis on efficiently allocating resources to meet fluctuating demands and optimize performance metrics. In the following, we provide an in-depth analysis of DRL-based resource provisioning strategies, focusing on their deployment within cloud and edge computing environments, while examining their capacity to improve system efficiency, reduce costs, and enhance QoS. Table IV offers a summary of these approaches.

1) *Value-Based DRL for Resource Provisioning*: Value-based DRL approaches are well-suited for optimizing resource provisioning by estimating action values to guide decision-making. DQN has been widely applied to address dynamic resource provisioning challenges, particularly in cloud and containerized environments. For example, the work [155]

TABLE IV: Summary of typical works on DRL-based resource provisioning

Reference	Method	Others	Environment	Makespan	Cost	Optimization objective	
						Energy consumption	Others
[155]	DQN	-	Cloud computing	-	✓	-	SLA requirements
[156]	DQN	-	Cloud computing	✓	✓	✓	SLA requirements
[157]	DDQN	-	Cloud computing	-	✓	-	Spectrum consumption
[158]	DQN	-	Edge computing	-	✓	-	Resource utilization
[159]	DQN	-	Edge-Cloud collaboration	-	-	-	Initial delay; blocking probability
[160]	DQN	-	Edge computing	✓	✓	✓	Resource utilization
[161]	DQN	-	Edge computing	-	-	✓	-
[162]	A3C	-	Cloud computing	-	-	✓	QoS
[163]	PPO	-	Cloud computing	-	✓	-	Resource utilization
[164]	DDPG	-	Cloud computing	-	-	-	Blocking probability
[165]	Actor-Critic	-	Edge computing	-	✓	-	Execution delay
[166]	A3C	-	Edge-Cloud collaboration	✓	✓	-	QoS
[167]	PPO	-	Edge computing	-	✓	✓	Resource utilization
[168]	DDPG	-	Edge computing	-	✓	-	-
[169]	DDPG	-	Edge computing	-	-	-	Resource utilization
[170]	MADRL	-	Cloud computing	✓	✓	✓	Resource utilization; load balancing
[79]	MADRL	QR-DQN	Cloud computing	-	-	✓	-
[171]	MADRL	CRO	Cloud computing	-	-	✓	-
[172]	MARL	-	Edge-Cloud collaboration	-	-	-	System latency
[173]	MADDPG	-	Edge computing	-	-	-	Training time cost
[174]	MAA2C	-	Edge computing	-	✓	-	Request delay
[175]	SARSA	GA	Cloud computing	✓	-	-	Resource utilization; load balancing
[176]	MADRL	Generative AI	Cloud computing	-	✓	✓	Execution latency
[177]	Q-learning	LSTM and LA	Edge computing	✓	-	✓	Offloading
[178]	-	POKTR	Edge computing	-	-	-	TUR
[179]	DQN	MCTS	Edge computing	-	✓	-	QoS
[180]	RL	MCTS	Edge computing	-	-	✓	Service latency

introduces an elastic resource provisioning method that integrates DQN to achieve efficient horizontal scaling for cloud services. By dynamically adjusting the number of VMs in response to fluctuating user demands, this approach minimizes resource wastage and reduces SLA violations, thus improving resource utilization and reducing overall costs. Similarly, the work [156] applies DQN to optimize VM configurations based on workload characteristics, aiming to minimize energy consumption, operational costs, and task response time, while simultaneously improving QoS and resource utilization. However, the inherent overestimation bias in DQN can limit its efficacy in resource provisioning. To address this, the work [157] proposes a DDQN-based solution for resource provisioning in inter-datacenter elastic optical networks. This framework optimizes the deployment and reuse of Virtual Network Functions (VNFs), enabling more efficient management of IT and spectrum resources to meet the demands of VNF service chaining. By decomposing complex VNF service chains (VNF-SCs) into manageable segments and applying an encoding scheme for standardized input lengths, DDQN facilitates efficient processing of diverse provisioning tasks.

In edge computing environments, where resource demands are often dynamic and heterogeneous, valued-based DRL methods have also been widely adopted to address unique provisioning challenges. For example, the work [158] introduces a DQN-driven framework designed to tackle resource provisioning challenges in MEC for 6G networks, enabling efficient, adaptive resource scaling and optimal service placement for IoE services. Similarly, the work [159] proposes a DQN-based solution for resource provisioning in cloud-edge environments connected via elastic optical networks (EONs). By optimizing both time and spectrum resources, this framework reduces network fragmentation and blocking probability, dynamically adapting to deadline-sensitive demands and improving re-

source utilization. Furthermore, the work [160] presents a self-learning DQN-based method for proactive resource and service provisioning in edge computing, targeting reductions in response time, cost, and energy consumption while improving overall resource utilization. To address the overestimation bias in conventional DQN, the work [161] introduces a DRL method called D3QN-PER, which combines D3QN with Prioritized Experience Replay for Service Function Chaining Allocation (SFCA) in fog computing. This enhanced method adapts to dynamic network conditions while improving energy efficiency.

2) *Policy-Based DRL for Resource Provisioning*: Policy-based DRL techniques have emerged as powerful tools for determining optimal resource provisioning decisions, enabling more efficient and adaptive management of resources in dynamic environments. In the context of cloud computing, several studies have explored the application of policy-based DRL methods to enhance resource allocation and management. For instance, the work [162] proposes an adaptive resource provisioning framework based on the A3C algorithm. This framework utilizes the policy-based DRL to handle dynamic system states and heterogeneous user demands, improving QoS and energy efficiency in cloud datacenters. Moreover, the work [163] proposes an automated cloud resource provisioning method based on PPO to address the issue of heterogeneous resource management. Specifically, this framework dynamically adjusts the number and types of VMs in a cloud environment, with the goal of enhancing system cost-efficiency and resource utilization by automatically adjusting resource provisioning. Furthermore, the work [164] introduced a DDPG-based service framework for resource provisioning in datacenter interconnections (DCIs), concentrating on virtual network slicing. This framework utilizes DDPG to assist in pricing and advertising substrate resources across non-

overlapping subgraphs of a DCI, allowing tenants to compute their own virtual network embedding (VNE) schemes independently, which reduces resource conflicts and computation time. This tenant-driven approach optimizes resource utilization and cost-effectiveness while improving scalability and efficiency compared to traditional centralized methods.

Similarly, in edge and fog computing environments, policy-based DRL strategies have been increasingly employed to address specific challenges related to resource provisioning and service optimization. For instance, the work [165] introduces the DRL-Dispatcher, which uses an AC algorithm to optimize task scheduling and resource provisioning for IoT requests, focusing on minimizing execution delays and operational costs. To address the complexities of resource provisioning and service migration in a heterogeneous cloud-edge environment supporting IoT applications, this work [166] utilizes A3C algorithm for trusted and dynamic Service Function Chain (SFC) orchestration. The approach aims to minimize orchestration costs and improve QoS, optimizing resource provisioning in both edge and cloud layers to meet the demands of high-mobility IoT networks. Moreover, PPO has been widely applied to address resource provisioning challenges in edge computing environments. The work [167] addresses fog computing challenges by utilizing the PPO algorithm to overcome issues such as uneven resource provisioning, suboptimal QoS, and low network bandwidth utilization. This approach aims to reducing latency, minimize deployment costs, and improve resource utilization. Furthermore, DDPG has been employed to tackle similar challenges in resource provisioning. In mobile edge computing, the work [168] proposes a DDPG-based method for computation offloading, which predicts optimal resource provisioning actions to reduce overall system costs. Additionally, the work [169] explores dynamic edge server reservation for connected vehicles in edge computing. This framework employs a DDPG algorithm enhanced with ConvLSTM and an action amender to capture spatio-temporal correlations in resource demand. By adapting server reservation decisions to real-time workload observations, this approach effectively addresses fluctuating resource requirements and minimizes provisioning inefficiencies.

3) *Multi-Agent DRL for Resource Provisioning*: MADRL has emerged as a powerful approach for optimizing resource provisioning, particularly in distributed systems where coordination among multiple agents is essential. Recent advancements in cloud computing have demonstrated the effectiveness of MADRL in addressing resource allocation and energy efficiency, particularly in tackling the challenges of dynamic provisioning and large-scale optimization. For example, the work [170] proposes a dynamic resource provisioning framework in cloud computing, focused on improving load balancing and service brokering through a MADRL approach. Specifically, the MADRL model anticipates user demand to prioritize resource allocation to VMs, thereby reducing response time and balancing loads across distributed nodes. Additionally, the framework introduces the dynamic optimal load-aware service broker strategy to optimize task scheduling among cloud brokers, aiming to minimize costs, improve energy efficiency, and meet QoS requirements. Similarly, the

work [79] employs a Quantile Regression-Deep Q Network (QR-DQN) algorithm within a multi-agent system to manage task allocation and resource provisioning, with a focus on reducing energy consumption while maintaining QoS in large-scale data centers. Furthermore, the work [171] presents a hybrid approach that combines the Coral Reef Optimization (CRO) algorithm with a Multi-Agent DQN (MDQ-CR) to enhance energy-aware resource provisioning using DVFS technology in cloud data centers. This two-phase model employs CRO for initial resource allocation and then utilizes a Multi-Agent Deep Q-Network for long-term resource management, addressing high energy demands by avoiding local optima and emphasizing long-term optimization.

As edge computing systems face increasingly complex challenges, MADRL has become a crucial method for addressing resource provisioning in distributed, and heterogeneous environments. For example, the work [172] addresses resource provisioning challenge in a distributed multi-cloud, MEC network by leveraging a MADRL approach. It considers a three-layer architecture involving cloud centers (CCs), MEC servers, and edge devices (EDs), where tasks are distributed among independent CCs that rely on MEC servers and EDs for data processing to minimize latency. The proposed solution optimizes task offloading and resource provisioning in a decentralized manner by allowing each CC to predict and adapt to the resource usage of other CCs, thus reducing system latency and enhancing resource utilization across heterogeneous devices in real time. Additionally, the work [173] introduces a MADDPG algorithm for resource provisioning in edge network slicing, aiming to balance latency and energy consumption. The proposed framework addresses the limitations of static slicing and instantaneous reward maximization by implementing a novel incremental learning scheme, which adapts the algorithm to dynamic changes in the number of slices without retraining from scratch. Furthermore, the work [174] proposes a MADRL approach for dynamic and efficient resource provisioning in 5G end-to-end network slicing, integrating MEC with network function virtualization (NFV) to support diverse user equipment (UE) demands in densely populated areas like airports and train stations. This strategy aims to maximize service provider profitability, uphold SLAs, and reduce operational costs by effectively managing network slices across distributed cloudlets in MEC environments.

4) *Advanced DRL for Resource Provisioning*: Advanced DRL techniques have expanded traditional DRL methods by integrating complementary approaches such as QRL [176] and LSTM-augmented DRL [177], offering innovative solutions to complex resource provisioning challenges. In the context of cloud computing, these techniques have enabled the development of sophisticated strategies to enhance resource efficiency. For example, the work [175] introduces a hybrid method that combines parallel SARSA RL agents with GA to improve resource provisioning. This approach focuses on reducing makespan and improving resource utilization within cloud environments. Furthermore, the work [176] presents an intelligent resource provisioning method designed for QRL, particularly focusing on generative AI applications across cloud, edge and mobile nodes. Specifically, this framework uti-

TABLE V: Summary of typical works on DRL-based resource scheduling

Reference	Method	Others	Environment	Makespan	Cost	Energy consumption	Optimization objective	Others
[52]	DQN	-	Cloud computing	✓	-	✓	-	-
[181]	DQN	-	Cloud computing	✓	-	-	-	Resource utilization
[182]	DQN	-	Cloud computing	-	✓	-	-	QoS
[183]	DDQN	-	Cloud computing	✓	✓	-	-	System throughput
[54]	DQN	-	Edge computing	-	-	-	-	Overhead and latency
[184]	DQN	-	Edge computing	✓	-	-	-	Network latency
[185]	DQN	-	Edge computing	-	-	-	-	Delay and computational cost
[186]	DQN	-	Edge computing	-	-	-	-	Network delay
[187]	DQN	-	Edge computing	✓	✓	-	-	-
[53]	DDQN	-	Edge computing	✓	✓	✓	-	Bandwidth cost
[188]	DDQN	-	Edge computing	-	-	-	-	Resource utilization; cost of bandwidth
[189]	PDQN	-	Edge computing	✓	-	-	-	Latency
[190]	REINFORCE	-	Cloud computing	✓	-	-	-	-
[191]	Actor-Critic	-	Cloud computing	✓	-	-	-	Throughput; resource utilization
[192]	A2C	-	Cloud computing	-	-	-	-	Job latency
[?]]	PPO	-	Cloud computing	-	✓	-	-	SLA requirements
[193]	A3C	-	Edge-Cloud collaboration	-	-	-	-	Delay; task drop rate
[194]	DDPG	-	Edge computing	-	✓	-	-	Task offloading
[195]	TD3	-	Edge computing	-	✓	-	-	Task offloading; processing delay
[196]	MADRL	Q-learning	Cloud computing	-	-	✓	-	Fault tolerance; workload balancing
[55]	MADRL	-	Cloud computing	-	-	-	-	Resource utilization
[197]	MADRL	-	Cloud computing	✓	-	✓	-	Resource utilization
[198]	MADRL	-	Edge computing	-	✓	✓	-	Computation offloading
[199]	MADRL	-	Edge computing	✓	✓	-	-	Resource usage
[56]	MADRL	DDQN	Edge computing	-	-	-	-	Delay; bandwidth
[103]	MADRL	-	Edge computing	✓	-	-	-	-
[200]	REINFORCE	Imitation learning	Cloud computing	✓	-	-	-	-
[201]	DQN	Simulated annealing	Cloud computing	-	✓	-	-	-
[202]	DRL variant	NESRL	Cloud computing	-	-	-	-	Resource utilization balance
[203]	DRL variant	2r-SAE; ASA; 2pER	Edge computing	✓	-	-	-	-
[204]	DQN	GA	Edge computing	✓	-	-	-	-
[205]	QRL	Grover search	Edge computing	-	-	✓	-	-
[206]	LQ-DRL	-	Edge computing	-	-	✓	-	QoS

lizes QNNs to represent policies, allowing the RL agents to optimize resource provisioning by taking advantage of quantum-specific phenomena, such as superposition and entanglement. By utilizing QRL, the framework achieves faster convergence rates, reduced decision-making latency, and enhanced stability, effectively improving resource management.

The development of advanced DRL techniques has significantly contributed to addressing the unique challenges posed by distributed architectures in edge computing environments. For instance, the work [177] introduces a hybrid method that integrates learning automata (LA), LSTM, and RL to facilitate dynamic resource provisioning decisions. Specifically, LA optimizes action selection by adjusting probabilities based on past outcomes and feedback from the environment, ultimately choosing the action with the highest probability. LSTM models analyze historical data to forecast future request volumes, enhancing the ability of RL agent to make timely and effective resource provisioning decisions. Additionally, the work [178] proposes an intelligent VNF configuration framework for the Cloud of Things. This method enhances the AC model by modifying the the loss function of the agent, which enables stable, monotonic improvements and accelerates convergence, essential for achieving reliable and efficient resource provisioning. To support sustainable resource provisioning for VNFs across multi-network operators, the work [179] combines DQN with MCTS. This method filters subcarrier(s)-end-user associations early in the decision process, significantly reducing the computational complexity of DQN for resource provisioning and expediting decision-making. By integrating MCTS, the approach optimizes resource provisioning to meet the computational demands of VNFs in edge networks. Similarly, the work [180] introduces

the iRAF, a resource provisioning model specifically tailored for MEC in IoT environments. By utilizing multitask DRL combined with MCTS, iRAF dynamically allocates resources across edge servers to manage latency and energy requirements in IoT applications. Adapting to real-time changes in network conditions, this framework enhances performance, optimizing resource provisioning in response to complex, fluctuating demands.

B. DRL-Based Resource Scheduling

Resource scheduling is a critical aspect of optimizing the performance and efficiency of computing systems, encompassing both cloud and edge environments. Subsequently, we provide a comprehensive overview of various DRL-based resource scheduling methods, detailing their specific applications and contributions to resource management within cloud and edge computing settings. These studies are systematically categorized in Table V, offering an in-depth summary of recent advancements in this field.

1) *Value-Based DRL for Resource Scheduling*: Value-based DRL methods concentrate on assessing the value of actions to optimize decision-making, providing a powerful approach for tackling resource scheduling challenges in dynamic and distributed computing environments [52], [181], [183]. DQN, as one of the foundational methods in value-based DRL, has been widely utilized to address resource scheduling problems in cloud computing environments. For instance, the work [52] develops a joint VM resource scheduling and power management framework for cloud computing systems, utilizing DQN to efficiently handle high-dimensional state and action spaces. The framework aims to minimize power consumption while ensuring acceptable performance levels. Furthermore,

the work [181] proposed a hybrid anomaly-aware DQN-based resource scaling method for dynamic resource allocation in cloud environments. DQN is employed to ensure efficient resource utilization and minimize makespan, while the approach also integrates anomaly detection to enhance decision-making stability in response to anomalous states. Additionally, the work [182] introduces a resource scheduling approach for cloud-based software services. This method employs a DQN-based prediction model trained using workload-time windows to anticipate management operations, specifically adjusting the number of VMs in response to varying system states. However, the inherent overestimation bias in traditional DQN can hinder learning stability. To address this, the work [183] incorporates DDQN into the proposed Deep Elastic Resource Provisioning (DERP) method. This approach aims to enhance elasticity in large-scale computing clusters by optimizing metrics such as makespan, resource cost, and system throughput.

Beyond cloud computing, value-based DRL approaches have also demonstrated its utility in edge computing resource scheduling. For example, the work [54] addresses the complexities of computing task offloading in MEC for the Internet of Vehicles (IoV) by developing a DQN-based resource scheduling scheme. The framework considers service node computational capabilities and vehicle velocity to achieve minimized system overhead and latency. Likewise, the work [184] designs a DQN-driven, cloud-edge cooperative strategy for content delivery within the IoV framework, effectively reducing network latency by optimizing caching and routing decisions. This approach leverages historical request patterns and real-time network states, enabling adaptive content management across network nodes. Furthermore, the work [185] introduces a DQN-based computing offloading resource scheduling strategy to address challenges such as increased latency, energy consumption, and reduced service quality in vehicular networks. Moreover, the work [186] proposes a DQN-based resource scheduling scheme to enhance content distribution in a layered fog radio access network (FRAN). It addresses the challenge of low-latency content transmission by formulating an optimal resource allocation problem as a minimal delay model, utilizing in-network caching to aggregate content requests. Furthermore, the work [187] proposes an action-constrained DQN approach for secure computing resource scheduling in serverless cloud-edge computing environments, with the objective of reducing overall system costs and makespan, while maintaining security in resource utilization.

Given the limitations of traditional DQN in managing the complexities of edge computing, several DQN variants have been proposed to improve performance in these scenarios. For instance, the work [53] employs DDQN to address the challenge of insufficient processing capabilities in wireless devices. By introducing a model for partial computation offloading and resource scheduling in mobile edge computing, this approach aims to optimize the system's weighted sum cost, computation delays, power consumption, and bandwidth utilization. Likewise, the work [188] proposes an edge-cloud resource scheduling framework leveraging DDQN to dynamically assess both resource utilization and network conditions, thereby enabling optimal offloading decisions that minimize

delays and fulfill computational and communication requirements. Another work [189] addresses the unique challenges of computation-intensive task processing on edge devices by proposing a method based on Parameterized DQN (PDQN) for service placement and resource scheduling. In this approach, service placement involves discrete decisions, whereas resource scheduling demands continuous decisions. PDQN handles this mixed action space by defining a deterministic function that maps states to the continuous parameters associated with each discrete action.

2) *Policy-Based DRL for Resource Scheduling*: Policy-based DRL methods directly learn the policy that maps states to actions, demonstrating significant advantages in optimizing resource scheduling across diverse environments. Techniques such as REINFORCE [190], AC [191], A2C [192] and PPO [207] have been explored in cloud computing for their efficiency and adaptability. For instance, the work [190] introduces the DeepRM method, which represents the scheduling policy through a DNN and trains it using the REINFORCE algorithm to minimize average job slowdown. Furthermore, the work [191] proposes an adaptive resource scheduling method based on AC to address the load balancing optimization problem in cloud data centers. According to the strategy learned by the actor, appropriate VMs are allocated to tasks, thereby optimizing makespan, response time, resource utilization, and throughput. Compared to REINFORCE, the AC method reduces variance, leading to faster convergence and improved stability. Moreover, the work [192] applies the A2C method to resource scheduling in data centers, dynamically adjusting task assignments based on current resource utilization to minimize task latency. By incorporating the advantage function, A2C reduces variance in policy gradient estimates, resulting in more stable learning and less noisy updates. PPO has also been applied to enhance sample efficiency in resource scheduling. The work [207] presents a novel PPO-based approach for automated resource scheduling in heterogeneous cloud environments. This method enables the PPO agent to interact with a data center environment containing large, medium, and small VMs, learning policies to maximize resource cost-effectiveness and meet SLA targets.

As edge computing gains prominence, policy-based DRL has also been adapted to meet the distinct demands of edge and IoT networks. For instance, in the field of blockchain, the work [193] addresses security challenges in edge-centric computing by utilizing a blockchain-based framework that integrates the A3C algorithm for resource scheduling. This approach enhances trust in the system while minimizing delays and task drop rates, providing a secure and efficient solution for resource scheduling in decentralized networks. In the realm of VEC, the work [194] formulates a two-stage Stackelberg game to incentivize computation offloading by VEC servers, using a DDPG-based resource scheduling scheme. The DDPG algorithm effectively manages continuous action spaces, allowing both vehicles and servers to optimize their resource utilization and maximize profits. Furthermore, the work [195] proposes a joint computation offloading and resource scheduling framework in IoV. This method leverages the TD3 algorithm, which offers improved stability and faster

convergence in handling continuous action spaces, thereby reducing system costs while maintaining high performance.

3) *Multi-Agent DRL for Resource Scheduling*: MADRL enables multiple agents to learn and interact within a shared environment, making it particularly effective for tackling the complexities of distributed systems. Its application in cloud computing has opened new avenues for improving efficiency and adaptability in resource scheduling. For example, the work [196] combines multi-agent framework with Q-learning by treating VMs as agents that dynamically adjust their states to meet requirements for energy consumption, fault tolerance, and load balancing, thereby optimizing overall system performance. In a similar vein, the work [55] explores the configuration of a high-performance AI computing environment using advanced technologies, such as Nvidia DGX servers, to leverage MADRL for regional resource utilization optimization, resulting in more effective resource scheduling. Additionally, the work [197] focuses on energy efficiency in cloud environments by utilizing container-based virtualization, which offers higher efficiency than traditional VMs. This approach employs MADRL to dynamically schedule tasks and allocate resources, significantly reducing energy consumption while accounting for VM overheads and workload patterns.

As edge computing and IoT networks grow in complexity, MADRL has become a vital tool for optimizing resource scheduling in these domains. For instance, the work [198] addresses resource competition in IoT edge computing by treating users as independent learning agents within a dynamic and uncertain environment. This approach uses a multi-agent Q-learning algorithm, allowing users to make optimal offloading decisions independently, minimizing long-term system costs without needing to know the actions of others. Similarly, another work [199] applies DRL to industrial IoT (IIoT) systems, employing multi-agent systems for decentralized decision-making, where one agent is responsible for device resource scheduling and the other manages network resource scheduling. This dual-agent design enhances scalability and adaptability to changing numbers of network nodes, providing greater flexibility in resource scheduling. In the context of 5G networks, the work [56] proposes a MEC framework utilizing a decentralized MADRL algorithm to address the challenges of low latency and high reliability in resource scheduling. By considering task priorities and channel variability, this framework effectively minimizes long-term costs related to delay and bandwidth, allowing edge clouds to support machine learning tasks with optimized resource scheduling. Furthermore, the work [103] introduces a distributed scheduler for resource allocation in edge computing networks, leveraging a GNN-based MADRL paradigm to address the complexities of resource scheduling for machine learning tasks in edge environments. This framework includes a heterogeneous graph attention network to manage interactions among distributed agents, along with a task selection mechanism and conflict resolution strategies, enhancing scheduling performance in multi-task scenarios across distributed edge clusters. By aiming to minimize task completion time and improve resource utilization efficiency, this approach offers significant advancements in edge resource management.

4) *Advanced DRL for Resource Scheduling*: Advanced DRL methods have increasingly integrated diverse techniques to enhance learning efficiency and address complex resource scheduling challenges. By incorporating approaches such as imitation learning [200], evolutionary strategies [201], or quantum computing [206], these methods tackle sophisticated resource scheduling challenges, achieving superior optimization and adaptability. In cloud computing, several studies have explored such techniques. For instance, the work [200] integrates imitation learning with the REINFORCE algorithm to accelerate model training and convergence. Imitation learning mimics the behavior of heuristic algorithms, allowing the RL agent to leverage expert strategies from these heuristics instead of relying on random exploration, significantly speeding up the learning process. Alternatively, some studies have combined heuristic algorithms with DRL. For example, the work [201] addresses cloud service configuration and adaptive resource scheduling challenges by integrating a modified DQN algorithm with SA. SA enhances DRL by gradually lowering the exploration temperature, allowing the agent to explore a broader range of actions in the early training stages and then prioritize exploitation as training progresses. This temperature-controlled exploration mechanism enables the agent efficiently balance exploration and exploitation, thereby improving convergence speed and enabling more effective resource scheduling in dynamic cloud environments. Additionally, the work [202] introduces an advanced DRL-based method that integrates Natural Evolution Strategy (NES) into A3C to optimize resource scheduling across distributed data centers. By utilizing NES to approximate the gradient of the reward function, this method enhances training efficiency and maintains exploration diversity, resulting in a more balanced utilization of resources. This approach effectively addresses the challenges of multi-dimensional resource allocation, achieving improved performance in distributed cloud environments.

Some studies have proposed advanced DRL methods specifically for IoT and edge computing environments to address resource scheduling challenges. For instance, the work [203] introduces an advanced RL framework for online resource scheduling in large-scale MEC systems. This is enhanced by a related and regularized stacked auto-encoder (2r-SAE) for data compression, an adaptive SA approach for search efficiency, and a preserved and prioritized experience replay (2pER) mechanism to improve policy training. RL can also assist heuristic algorithms, for example, the work [204] combines DQN with GA to improve resource scheduling on edge cloud servers, minimize application execution time. In this approach, DQN generates the initial population for GA, reducing computational costs and improving optimization for the scheduling problem. Some studies further leverage QRL to boost resource scheduling performance. For example, the work [205] proposes an innovative QRL method for dynamic IoT networks that jointly optimizes resource scheduling, content caching, and computation offloading to maximize energy efficiency. Leveraging quantum computing techniques, including an improved Grover search algorithm, this method accelerates the training process and increases policy adaptability within high-dimensional continuous action

spaces, thereby significantly boosting overall effectiveness. Similarly, the work [206] develops a layerwise QRL to address continuous large-space and time-series challenges for resource scheduling. By utilizing quantum embeddings, this framework focuses on UAV trajectory planning and power allocation, aiming to optimize energy consumption while ensuring QoS.

VI. FUTURE DEVELOPMENT TRENDS

Building on the review of existing works in DRL-based job scheduling and resource management, this section discusses the challenges of applying DRL in cloud computing and provides insights into potential future directions.

A. Privacy and Security for Job Scheduling

As cloud computing continues to expand, job scheduling faces significant challenges in addressing privacy and security concerns, particularly in scenarios involving sensitive data processing and cross-platform collaboration. Privacy protection, especially in shared cloud environments, requires the implementation of strict data and task privacy constraints during scheduling. To address this, one promising direction involves leveraging hybrid cloud architectures, which enable resource isolation and elastic management to support dynamic partitioning and cross-data-center scheduling [208], [209]. Another critical avenue is the integration of privacy protection directly into the objective function of scheduling algorithms, allowing for a balanced optimization of privacy, performance, and cost. [210] Security remains a parallel concern, with the confidentiality and integrity of sensitive data posing significant challenges. Current practices, such as employing encryption algorithms like RC4 for confidentiality [211] and SHA-1 for integrity [212], offer foundational safeguards for intermediate and stored data. Future research will focus on refining these security measures, particularly by integrating them into the scheduling process. This includes addressing issues such as user authentication, secure resource allocation, and scalable security protocols to ensure robust protection for sensitive data in increasingly complex and heterogeneous cloud environments.

B. Resource Management over Multi-Tier Networks

Modern computing environments increasingly rely on multi-tier networks to address growing complexities and diverse application demands. These networks, comprising edge, fog, mist, and dew computing layers, pose significant challenges for resource management due to their heterogeneous resource characteristics and varied latency requirements. Addressing these challenges requires advanced frameworks and innovative technologies tailored to the unique needs of multi-tier architectures. One promising direction is the integration of DT technology into resource management [213]. DT technology provides synchronized virtual representations of physical entities, enabling real-time monitoring, predictive analysis, and proactive management of resource utilization and user demands. Another significant future trend is the adoption of multi-tier multi-domain network slicing [214], [215]. This approach facilitates resource aggregation across multiple infrastructure providers,

allowing for more effective distribution of resources across domains and the implementation of tailored strategies that meet the specific requirements of diverse applications.

C. Handling Large-Scale Decision-Making in Job Scheduling and Resource Management

In large-scale job scheduling and resource management scenarios, such as multi-cloud data centers, the inherent complexity of these environments introduces significant challenges. For single-agent systems, the primary obstacle lies in the overwhelming state and action spaces, commonly referred to as the “curse of dimensionality”. The high dimensionality of these spaces increases the computational complexity and hampers the learning efficiency of DRL algorithms. Hierarchical reinforcement learning emerges as a promising solution to address these challenges by decomposing job scheduling or resource management into multiple decision-making layers [216], [217]. High-level agents manage coarse-grained decisions, while low-level agents handle finer-grained decisions. This layered approach simplifies individual agent tasks, improving both learning efficiency and scalability. In multi-agent systems, coordinating hundreds or thousands of agents introduces further challenges. Traditional centralized approaches face significant limitations, including high communication overhead, and latency, alongside scalability issues as system size increases. To overcome these limitations, decentralized model-based policy optimization frameworks have emerged as a promising research direction [218]. These frameworks enable agents to interact primarily with immediate neighbors, significantly reducing communication costs while maintaining effective coordination. By leveraging localized interactions and model-based learning, decentralized approaches enhance scalability and efficiency. This paradigm allows agents to achieve robust decision-making without relying on extensive global communication, making it particularly suitable for large-scale and complex environments.

D. Incorporating Large Language Models (LLMs) for Job Scheduling and Resource Management

Job scheduling and resource management in dynamic environments face notable challenges due to their complexity and variability. Traditional DRL approaches, while effective in well-defined scenarios, often struggle to adapt to unforeseen situations, relying heavily on training tailored to specific environments, such as fixed-scale resource pools or predetermined user request patterns. This lack of flexibility limits their applicability in diverse and dynamic contexts, while frequent retraining to accommodate changing conditions adds significant computational overhead. To address these challenges, the integration of LLMs represents a transformative research direction. Unlike traditional DRL methods, LLMs possess a remarkable ability to generalize across diverse scenarios, enabling them to adapt to previously unseen environments without the need for exhaustive retraining [219]. This adaptability is especially valuable in dynamic systems, where LLMs can interpret changing workloads and resource states in real time, providing more flexible and efficient

solutions. Furthermore, their capability for knowledge transfer across different contexts accelerates learning process, reduces deployment overhead, and supports robust decision-making in heterogeneous environments. By leveraging these strengths, LLMs open new avenues for advancing job scheduling and resource management in complex, ever-changing settings.

VII. CONCLUSION

This paper presents a comprehensive analysis of advancements in deep reinforcement learning (DRL) methodologies for job scheduling and resource management in cloud computing, with a focus on reviewing existing works categorized by the specific DRL algorithms used. We began by outlining the modeling of these optimization problems using Markov Decision Processes (MDPs) and demonstrated how DRL can effectively address them. Subsequently, we provided an overview of current DRL algorithms and systematically reviewed their applications in job scheduling, including task scheduling and workflow scheduling, as well as resource management, focusing on resource provisioning and scheduling. The reviewed works are categorized based on the DRL algorithms employed, offering a clear framework for understanding their implementation and impact.

In addition to analyzing existing methodologies, we provide insights to guide future research and practical advancements, paving the way for more efficient and adaptive DRL-driven solutions. This review underscores the pivotal role of DRL in tackling the complex and dynamic challenges of job scheduling and resource management. As cloud computing environments continue to evolve, we expect that DRL-based approaches will become increasingly proficient at handling the unpredictable demands of workload and resource allocation, driving significant advancements in these domains.

REFERENCES

- [1] D. C. Marinescu, *Cloud Computing: Theory and Practice*. Morgan Kaufmann, 2022.
- [2] L. Cheng, B. F. van Dongen, and W. M. van der Aalst, "Scalable discovery of hybrid process models in a cloud computing environment," *IEEE Transactions on Services Computing*, vol. 13, no. 2, pp. 368–380, 2020.
- [3] Y. Mao, W. Yan, Y. Song, Y. Zeng, M. Chen, L. Cheng, and Q. Liu, "Differentiate quality of experience scheduling for deep learning inferences with docker containers in the cloud," *IEEE Transactions on Cloud Computing*, vol. 11, no. 2, pp. 1667–1677, 2023.
- [4] S. Duan, D. Wang, J. Ren, F. Lyu, Y. Zhang, H. Wu, and X. Shen, "Distributed artificial intelligence empowered by end-edge-cloud computing: A survey," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 591–624, 2022.
- [5] Q. Luo, S. Hu, C. Li, G. Li, and W. Shi, "Resource scheduling in edge computing: A survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 4, pp. 2131–2165, 2021.
- [6] J. Lu, J. Yang, S. Li, Y. Li, W. Jiang, J. Dai, and J. Hu, "A2c-drl: Dynamic scheduling for stochastic edge-cloud environments using a2c and deep reinforcement learning," *IEEE Internet of Things Journal*, 2024.
- [7] F. Cheng, Y. Huang, B. Tanpure, P. Sawalani, L. Cheng, and C. Liu, "Cost-aware job scheduling for cloud instances using deep reinforcement learning," *Cluster Computing*, pp. 1–13, 2022.
- [8] A. Chatterjee, P. Dubey, and A. Nigam, "Dynamic On-Demand Machine Provisioning and Continuous Resource Management," in *2023 International Conference on Computing, Communication, and Intelligent Systems*, 2023, pp. 1010–1015.
- [9] Y. Lei and S. Jasin, "Real-time dynamic pricing for revenue management with reusable resources, advance reservation, and deterministic service time requirements," *Operations Research*, vol. 68, no. 3, pp. 676–685, 2020.
- [10] P. Murthy, A. Mehra, and L. Mishra, "Resource Allocation for Generative AI Workloads: Advanced Cloud Resource Management Strategies for Optimized Model Performance," *Iconic Research And Engineering Journals*, vol. 6, p. 12, 2023.
- [11] J. Yan, Y. Huang, A. Gupta, A. Gupta, C. Liu, J. Li, and L. Cheng, "Energy-aware systems for real-time job scheduling in cloud data centers: A deep reinforcement learning approach," *Computers and Electrical Engineering*, vol. 99, p. 107688, 2022.
- [12] J. Liu, H. Shen, H. Chi, H. S. Narman, Y. Yang, L. Cheng, and W. Chung, "A low-cost multi-failure resilient replication scheme for high-data availability in cloud storage," *IEEE/ACM Transactions on Networking*, vol. 29, no. 4, pp. 1436–1451, 2021.
- [13] S. Yang, F. Li, S. Trajanovski, R. Yahyapour, and X. Fu, "Recent advances of resource allocation in network function virtualization," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 2, pp. 295–314, 2020.
- [14] Y. Xu, G. Gui, H. Gacanin, and F. Adachi, "A survey on resource allocation for 5G heterogeneous networks: Current research, future trends, and challenges," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 2, pp. 668–695, 2021.
- [15] E. H. Houssein, A. G. Gad, Y. M. Wazery, and P. N. Suganthan, "Task scheduling in cloud computing based on meta-heuristics: review, taxonomy, open challenges, and future trends," *Swarm and Evolutionary Computation*, vol. 62, p. 100841, 2021.
- [16] A. Pradhan, S. K. Bisoy, and A. Das, "A survey on PSO based meta-heuristic scheduling mechanism in cloud computing environment," *Journal of King Saud University-Computer and Information Sciences*, vol. 34, no. 8, pp. 4888–4901, 2022.
- [17] X. Liu and R. Buyya, "Resource management and scheduling in distributed stream processing systems: a taxonomy, review, and future directions," *ACM Computing Surveys*, vol. 53, no. 3, pp. 1–41, 2020.
- [18] B. Jamil, H. Ijaz, M. Shojafar, K. Munir, and R. Buyya, "Resource allocation and task scheduling in fog computing and internet of everything environments: A taxonomy, review, and future directions," *ACM Computing Surveys*, vol. 54, no. 11s, pp. 1–38, 2022.
- [19] M. Afrin, J. Jin, A. Rahman, A. Rahman, J. Wan, and E. Hossain, "Resource allocation and service provisioning in multi-agent cloud robotics: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 2, pp. 842–870, 2021.
- [20] G. Zhou, W. Tian, R. Buyya, R. Xue, and L. Song, "Deep reinforcement learning-based methods for resource scheduling in cloud computing: A review and future directions," *Artificial Intelligence Review*, vol. 57, no. 5, p. 124, 2024.
- [21] Z. Jalali Khalil Abadi, N. Mansouri, and M. M. Javidi, "Deep reinforcement learning-based scheduling in distributed systems: a critical review," *Knowledge and Information Systems*, pp. 1–74, 2024.
- [22] I. A. Mohialdeen, "Comparative study of scheduling algorithms in cloud computing environment," *Journal of Computer Science*, vol. 9, no. 2, pp. 252–263, 2013.
- [23] H. Chen, F. Wang, N. Helian, and G. Akanmu, "User-priority guided Min-Min scheduling algorithm for load balancing in cloud computing," in *2013 National Conference on Parallel Computing Technologies*, 2013, pp. 1–8.
- [24] R. NoorianTalouki, M. H. Shirvani, and H. Motameni, "A heuristic-based task scheduling algorithm for scientific workflows in heterogeneous cloud computing platforms," *Journal of King Saud University-Computer and Information Sciences*, vol. 34, no. 8, pp. 4902–4913, 2022.
- [25] F. Hoseiny, S. Azizi, M. Shojafar, F. Ahmadiazar, and R. Tafazolli, "PGA: a priority-aware genetic algorithm for task scheduling in heterogeneous fog-cloud computing," in *IEEE INFOCOM 2021-IEEE Conference on Computer Communications Workshops*, 2021, pp. 1–6.
- [26] X. Chen, L. Cheng, C. Liu, Q. Liu, J. Liu, Y. Mao, and J. Murphy, "A WOA-based optimization approach for task scheduling in cloud computing systems," *IEEE Systems Journal*, vol. 14, no. 3, pp. 3117–3128, 2020.
- [27] Z. Pooranian, M. Shojafar, P. G. V. Naranjo, L. Chiaraviglio, and M. Conti, "A Novel Distributed Fog-Based Networked Architecture to Preserve Energy in Fog Data Centers," in *2017 IEEE 14th International Conference on Mobile Ad Hoc and Sensor Systems*, 2017, pp. 604–609.
- [28] L. Ni, J. Zhang, C. Jiang, C. Yan, and K. Yu, "Resource Allocation Strategy in Fog Computing Based on Priced Timed Petri Nets," *IEEE Internet of Things Journal*, vol. 4, no. 5, pp. 1216–1228, 2017.

- [29] M. Yang, H. Ma, S. Wei, Y. Zeng, Y. Chen, and Y. Hu, "A Multi-Objective Task Scheduling Method for Fog Computing in Cyber-Physical-Social Services," *IEEE Access*, vol. 8, pp. 65 085–65 095, 2020.
- [30] N. Potu, C. Jatoth, and P. Parvataneni, "Optimizing resource scheduling based on extended particle swarm optimization in fog computing environments," *Concurrency and Computation: Practice and Experience*, vol. 33, no. 23, p. e6163, 2021.
- [31] L. Cheng, Y. Wang, F. Cheng, C. Liu, Z. Zhao, and Y. Wang, "A deep reinforcement learning-based preemptive approach for cost-aware cloud job scheduling," *IEEE Transactions on Sustainable Computing*, vol. 9, no. 3, pp. 422–432, 2024.
- [32] R. S. Sutton, "Reinforcement learning: An introduction," *A Bradford Book*, 2018.
- [33] Y. Li, "Deep reinforcement learning: An overview," *arXiv preprint arXiv:1701.07274*, 2017.
- [34] P. Ladosz, L. Weng, M. Kim, and H. Oh, "Exploration in deep reinforcement learning: A survey," *Information Fusion*, vol. 85, pp. 1–22, 2022.
- [35] H. Jiang, N. Bhujel, Z. Lin, K.-W. Wan, J. Li, S. Jayavelu, and X. Jiang, "Learning relation in crowd using gated graph convolutional networks for drl-based robot navigation," *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [36] S. Dai, S. Li, H. Tang, X. Ning, F. Fang, Y. Fu, Q. Wang, and L. Cheng, "Marp: A cooperative multiagent drl system for connected autonomous vehicle platooning," *IEEE Internet of Things Journal*, vol. 11, no. 20, pp. 32 454–32 463, 2024.
- [37] Q. Liu, L. Cheng, A. L. Jia, and C. Liu, "Deep reinforcement learning for communication flow control in wireless mesh networks," *IEEE Network*, vol. 35, no. 2, pp. 112–119, 2021.
- [38] Q. Liu, T. Xia, L. Cheng, M. Van Eijk, T. Ozcelebi, and Y. Mao, "Deep reinforcement learning for load-balancing aware network control in iot edge systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 6, pp. 1491–1502, 2022.
- [39] S. Smachat and K. Viriyapant, "Taxonomies of workflow scheduling problem and techniques in the cloud," *Future Generation Computer Systems*, vol. 52, pp. 1–12, 2015.
- [40] H. Liu, X. Zhou, K. Gao, and Y. Ju, "An integrated optimization method to task scheduling and vm placement for green datacenters," *Simulation Modelling Practice and Theory*, vol. 135, p. 102962, 2024.
- [41] H. He, Y. Gu, Q. Liu, H. Wu, and L. Cheng, "Job scheduling in hybrid clouds with privacy constraints: A deep reinforcement learning approach," *Concurrency and Computation: Practice and Experience*, p. e8307, 2024.
- [42] S. Song, Y.-K. Kwok, and K. Hwang, "Security-driven heuristics and a fast genetic algorithm for trusted grid job scheduling," in *19th IEEE International Parallel and Distributed Processing Symposium*, 2005.
- [43] L. Chen, S. Liu, B. Li, and B. Li, "Scheduling jobs across geo-distributed datacenters with max-min fairness," *IEEE Transactions on Network Science and Engineering*, vol. 6, no. 3, pp. 488–500, 2018.
- [44] A. Mazrekaj, I. Shabani, and B. Sejdiu, "Pricing schemes in cloud computing: an overview," *International Journal of Advanced Computer Science and Applications*, vol. 7, no. 2, pp. 80–86, 2016.
- [45] J. Guo, L. Bhuyan, R. Kumar, and S. Basu, "Qos aware job scheduling in a cluster-based web server for multimedia applications," in *19th IEEE International Parallel and Distributed Processing Symposium*, 2005.
- [46] S. Li, H. Jin, Y. Gao, Y. Wang, S. Dai, Y. Xu, and L. Cheng, "Approximate data mapping in refresh-free dram for energy-efficient computing in modern mobile systems," *Computer Communications*, vol. 216, pp. 151–158, 2024.
- [47] M. I. Khaleel, M. Safran, S. Alfarhood, and M. Zhu, "Energy-latency trade-off analysis for scientific workflow in cloud environments: the role of processor utilization ratio and mean grey wolf optimizer," *Engineering Science and Technology, an International Journal*, vol. 50, p. 101611, 2024.
- [48] K. Deng, K. Ren, M. Zhu, and J. Song, "A data and task co-scheduling algorithm for scientific cloud workflows," *IEEE Transactions on Cloud Computing*, vol. 8, no. 2, pp. 349–362, 2015.
- [49] H. Topcuoglu, S. Hariri, and M.-Y. Wu, "Performance-effective and low-complexity task scheduling for heterogeneous computing," *IEEE transactions on parallel and distributed systems*, vol. 13, no. 3, pp. 260–274, 2002.
- [50] Q. Wu, M. Zhou, and J. Wen, "Endpoint communication contention-aware cloud workflow scheduling," *IEEE Transactions on Automation Science and Engineering*, vol. 19, no. 2, pp. 1137–1150, 2022.
- [51] A. Mampage, S. Karunasekera, and R. Buyya, "A deep reinforcement learning based algorithm for time and cost optimized scaling of serverless applications," *arXiv preprint arXiv:2308.11209*, 2023.
- [52] N. Liu, Z. Li, J. Xu, Z. Xu, S. Lin, Q. Qiu, J. Tang, and Y. Wang, "A Hierarchical Framework of Cloud Resource Allocation and Power Management Using Deep Reinforcement Learning," in *2017 IEEE 37th International Conference on Distributed Computing Systems*, 2017, pp. 372–382.
- [53] H. C. Ke, H. Wang, H. W. Zhao, and W. J. Sun, "Deep reinforcement learning-based computation offloading and resource allocation in security-aware mobile edge computing," *Wireless Networks*, vol. 27, no. 5, pp. 3357–3373, 2021.
- [54] Y. Zhang, M. Zhang, C. Fan, F. Li, and B. Li, "Computing resource allocation scheme of IOV using deep reinforcement learning in edge computing environment," *EURASIP Journal on Advances in Signal Processing*, vol. 2021, no. 1, p. 33, 2021.
- [55] J. Narantuya, J.-S. Shin, S. Park, and J. Kim, "Multi-Agent Deep Reinforcement Learning-Based Resource Allocation in HPC/AI Converged Cluster," *Computers, Materials & Continua*, vol. 72, no. 3, pp. 4375–4395, 2022.
- [56] H. Ke, H. Wang, and H. Sun, "Multi-Agent Deep Reinforcement Learning-Based Partial Task Offloading and Resource Allocation in Edge Computing Environment," *Electronics*, vol. 11, no. 15, p. 2394, 2022.
- [57] Z. Lan, Z. Zeng, B. Hong, Z. Liu, and F. Ma, "Rcsearcher: Reaction center identification in retrosynthesis via deep Q-learning," *Pattern Recognition*, vol. 150, p. 110318, 2024.
- [58] X. Chen, Q. Yu, S. Dai, P. Sun, H. Tang, and L. Cheng, "Deep reinforcement learning for efficient iot data compression in smart railroad management," *IEEE Internet of Things Journal*, vol. 11, no. 15, pp. 25 494–25 504, 2024.
- [59] Y. Liu, A. Halev, and X. Liu, "Policy learning with constraints in model-free reinforcement learning: A survey," in *The 30th International Joint Conference on Artificial Intelligence*, 2021.
- [60] V. Konda and J. Tsitsiklis, "Actor-critic algorithms," *Advances in Neural Information Processing Systems*, vol. 12, 1999.
- [61] R. Mukhopadhyay, S. Bandyopadhyay, A. Sutradhar, and P. Chattopadhyay, "Performance analysis of deep q networks and advantage actor critic algorithms in designing reinforcement learning-based self-tuning pid controllers," in *2019 IEEE Bombay Section Signature Conference*. IEEE, 2019, pp. 1–6.
- [62] D. T. Gillespie, "Exact numerical simulation of the ornstein-uhlenbeck process and its integral," *Phys. Rev. E*, vol. 54, pp. 2084–2091, 1996.
- [63] B. Sarkar, A. Talati, A. Shih, and D. Sadigh, "Pantheonrl: A marl library for dynamic training interactions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 11, 2022, pp. 13 221–13 223.
- [64] P. Leroy, P. G. Morato, J. Pisane, A. Kolios, and D. Ernst, "Imp-marl: a suite of environments for large-scale infrastructure management planning via marl," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [65] L. Kraemer and B. Banerjee, "Multi-agent reinforcement learning as a rehearsal for decentralized planning," *Neurocomputing*, vol. 190, pp. 82–94, 2016.
- [66] G. Shen, L. Lei, X. Zhang, Z. Li, S. Cai, and L. Zhang, "Multi-uav cooperative search based on reinforcement learning with a digital twin driven training framework," *IEEE Transactions on Vehicular Technology*, vol. 72, pp. 8354–8368, 2023.
- [67] S. Li, Y. Wu, X. Cui, H. Dong, F. Fang, and S. Russell, "Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 4213–4220.
- [68] C. Daskalakis, D. J. Foster, and N. Golowich, "Independent policy gradient methods for competitive reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 5527–5540, 2020.
- [69] M. Ye, Q. Han, L. Ding, and S. Xu, "Distributed nash equilibrium seeking in games with partial decision information: A survey," *Proceedings of the IEEE*, vol. 111, pp. 140–157, 2023.
- [70] H. Hu, A. Lerer, A. Peysakhovich, and J. Foerster, "'other-pla' for zero-shot coordination," in *International Conference on Machine Learning*, 2020, pp. 4399–4410.
- [71] K. Zhang, S. Kakade, T. Basar, and L. Yang, "Model-based multi-agent rl in zero-sum markov games with near-optimal sample complexity," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1166–1178, 2020.

- [72] F. Zhang, J. Li, and Z. Li, "A td3-based multi-agent deep reinforcement learning method in mixed cooperation-competition environment," *Neurocomputing*, vol. 411, pp. 206–215, 2020.
- [73] A. Morvan, B. Villalonga, X. Mi, and et al., "Phase transitions in random circuit sampling," *Nature*, vol. 634, no. 8033, pp. 328–333, Oct. 2024.
- [74] M. Da Silva, C. Ryan-Anderson, J. Bello-Rivas, A. Chernoguzov, J. Dreiling, C. Foltz, J. Gaebler, T. Gatterman, D. Hayes, N. Hewitt et al., "Demonstration of logical qubits and repeated error correction with better-than-physical error rates," *arXiv preprint*, 2024.
- [75] S. Mangalampalli, G. R. Karri, M. Kumar, O. I. Khalaf, C. A. T. Romero, and G. A. Sahib, "Drlbtsa: Deep reinforcement learning based task-scheduling algorithm in cloud computing," *Multimedia Tools and Applications*, vol. 83, no. 3, pp. 8359–8387, 2024.
- [76] M. Cheng, J. Li, and S. Nazarian, "Drl-cloud: Deep reinforcement learning-based resource provisioning and task scheduling for cloud service providers," in *2018 23rd Asia and South Pacific design automation conference*. IEEE, 2018, pp. 129–134.
- [77] K. Kang, D. Ding, H. Xie, Q. Yin, and J. Zeng, "Adaptive drl-based task scheduling for energy-efficient cloud computing," *IEEE Transactions on Network and Service Management*, vol. 19, no. 4, pp. 4948–4961, 2021.
- [78] Y. Ran, H. Hu, Y. Wen, and X. Zhou, "Optimizing energy efficiency for data center via parameterized deep reinforcement learning," *IEEE Transactions on Services Computing*, vol. 16, no. 2, pp. 1310–1323, 2022.
- [79] T. Oudaa, H. Gharsellaoui, and S. B. Ahmed, "An agent-based model for resource provisioning and task scheduling in cloud computing using drl," *Procedia Computer Science*, vol. 192, pp. 3795–3804, 2021.
- [80] Y. Yang and H. Shen, "Deep reinforcement learning enhanced greedy optimization for online scheduling of batched tasks in cloud hpc systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 11, pp. 3003–3014, 2021.
- [81] Z. Tong, F. Ye, B. Liu, J. Cai, and J. Mei, "Ddqn-ts: A novel bi-objective intelligent scheduling algorithm in the cloud environment," *Neurocomputing*, vol. 455, pp. 419–430, 2021.
- [82] F. Han, N. Yu, J. Gong, Y. Ge, and X. Gao, "Task scheduling for mobile edge computing leveraging deep reinforcement learning," in *2023 42nd Chinese Control Conference*. IEEE, 2023, pp. 1921–1926.
- [83] H. Yuan, G. Tang, X. Li, D. Guo, L. Luo, and X. Luo, "Online dispatching and fair scheduling of edge computing tasks: A learning-based approach," *IEEE Internet of Things Journal*, vol. 8, no. 19, pp. 14985–14998, 2021.
- [84] L. Zeng, Q. Liu, S. Shen, and X. Liu, "Improved double deep q network-based task scheduling algorithm in edge computing for makespan optimization," *Tsinghua Science and Technology*, vol. 29, no. 3, pp. 806–817, 2023.
- [85] Q. Zhang, M. Lin, L. T. Yang, Z. Chen, S. U. Khan, and P. Li, "A double deep q-learning model for energy-efficient edge scheduling," *IEEE Transactions on Services Computing*, vol. 12, no. 5, pp. 739–749, 2018.
- [86] J. Jin and Y. Xu, "Optimal policy characterization enhanced proximal policy optimization for multitask scheduling in cloud computing," *IEEE Internet of Things Journal*, vol. 9, no. 9, pp. 6418–6433, 2021.
- [87] Y. Yang, C. He, B. Yin, Z. Wei, and B. Hong, "Cloud task scheduling based on proximal policy optimization algorithm for lowering energy consumption of data center," *KSII Transactions on Internet & Information Systems*, vol. 16, no. 6, 2022.
- [88] J. Zhao, M. A. Rodríguez, and R. Buyya, "A deep reinforcement learning approach to resource management in hybrid clouds harnessing renewable energy and task scheduling," in *2021 IEEE 14th International Conference on Cloud Computing*. IEEE, 2021, pp. 240–249.
- [89] L. Ran, X. Shi, and M. Shang, "Slas-aware online task scheduling based on deep reinforcement learning method in cloud environment," in *2019 IEEE 21st International Conference on High Performance Computing and Communications*, 2019, pp. 1518–1525.
- [90] Z. Zhao, X. Shi, and M. Shang, "Performance and cost-aware task scheduling via deep reinforcement learning in cloud environment," in *International Conference on Service-Oriented Computing*. Springer, 2022, pp. 600–615.
- [91] F. Zhang, L. Jiang, and J. Chen, "Etpam: An efficient task pre-assignment and migration algorithm in heterogeneous edge-cloud computing environments," in *2024 27th International Conference on Computer Supported Cooperative Work in Design*. IEEE, 2024, pp. 2400–2405.
- [92] S. Tuli, S. Ilager, K. Ramamohanarao, and R. Buyya, "Dynamic scheduling for stochastic edge-cloud computing environments using a3c learning and residual recurrent neural networks," *IEEE transactions on mobile computing*, vol. 21, no. 3, pp. 940–954, 2020.
- [93] F. Zhang, Z. Tang, J. Lou, and W. Jia, "Online joint scheduling of delay-sensitive and computation-oriented tasks in edge computing," in *2019 15th International Conference on Mobile Ad-Hoc and Sensor Networks*. IEEE, 2019, pp. 303–308.
- [94] Y. Chen, Y. Sun, C. Wang, and T. Taleb, "Dynamic task allocation and service migration in edge-cloud iot system based on deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 16742–16757, 2022.
- [95] H. A. Balla, C. G. Sheng, and W. Jing, "Reliability-aware: task scheduling in cloud computing using multi-agent reinforcement learning algorithm and neural fitted q," *Int. Arab J. Inf. Technol.*, vol. 18, no. 1, pp. 36–47, 2021.
- [96] S. Jung, W. J. Yun, M. Shin, J. Kim, and J.-H. Kim, "Orchestrated scheduling and multi-agent deep reinforcement learning for cloud-assisted multi-UAV charging systems," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 6, pp. 5362–5377, 2021.
- [97] M. I. Gergely, "Multi-agent deep reinforcement learning for collaborative task scheduling," in *ICAART (3)*, 2024, pp. 1076–1083.
- [98] Y. Zhang, R. Li, Y. Zhao, R. Li, Y. Wang, and Z. Zhou, "Multi-agent deep reinforcement learning for online request scheduling in edge cooperation networks," *Future Generation Computer Systems*, vol. 141, pp. 258–268, 2023.
- [99] Q. Tang, R. Xie, F. R. Yu, T. Chen, R. Zhang, T. Huang, and Y. Liu, "Distributed task scheduling in serverless edge computing networks for the internet of things: A learning approach," *IEEE Internet of Things Journal*, vol. 9, no. 20, pp. 19634–19648, 2022.
- [100] C. Xu, S. Liu, C. Zhang, Y. Huang, Z. Lu, and L. Yang, "Multi-agent reinforcement learning based distributed transmission in collaborative cloud-edge systems," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 2, pp. 1658–1672, 2021.
- [101] Z. Zhang, F. Zhang, Z. Xiong, K. Zhang, and D. Chen, "Lsia3cs: Deep reinforcement learning-based cloud-edge collaborative task scheduling in large-scale iiot," *IEEE Internet of Things Journal*, 2024.
- [102] Z. Li, J. Yu, X. Liu, and L. Peng, "Load balancing for task scheduling based on multi-agent reinforcement learning in cloud-edge-end collaborative environments," in *Proceedings of the 2024 8th International Conference on Machine Learning and Soft Computing*, 2024, pp. 94–100.
- [103] Y. Li, X. Zhang, T. Zeng, J. Duan, C. Wu, D. Wu, and X. Chen, "Task placement and resource allocation for edge machine learning: A gnn-based multi-agent reinforcement learning paradigm," *IEEE Transactions on Parallel and Distributed Systems*, 2023.
- [104] L. Niu, X. Chen, N. Zhang, Y. Zhu, R. Yin, C. Wu, and Y. Cao, "Multiagent meta-reinforcement learning for optimized task scheduling in heterogeneous edge computing systems," *IEEE Internet of Things Journal*, vol. 10, no. 12, pp. 10519–10531, 2023.
- [105] Y. Gong, H. Yao, J. Wang, L. Jiang, and F. R. Yu, "Multi-agent driven resource allocation and interference management for deep edge networks," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 2, pp. 2018–2030, 2021.
- [106] Y. Huang, L. Cheng, L. Xue, C. Liu, Y. Li, J. Li, and T. Ward, "Deep adversarial imitation reinforcement learning for qos-aware cloud job scheduling," *IEEE Systems Journal*, vol. 16, no. 3, pp. 4232–4242, 2022.
- [107] G. Rjoub, J. Bentahar, O. A. Wahab, and A. Bataineh, "Deep smart scheduling: A deep learning approach for automated big data scheduling over the cloud," in *2019 7th International Conference on Future Internet of Things and Cloud*. IEEE, 2019, pp. 189–196.
- [108] T. Li, S. Ying, Y. Zhao, and J. Shang, "Batch jobs load balancing scheduling in cloud computing using distributional reinforcement learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 35, no. 1, pp. 169–185, 2023.
- [109] X. Mao, G. Wu, M. Fan, Z. Cao, and W. Pedrycz, "DL-DRL: A double-level deep reinforcement learning approach for large-scale task scheduling of multi-UAV," *IEEE Transactions on Automation Science and Engineering*, 2024.
- [110] Z. Tang, W. Jia, X. Zhou, W. Yang, and Y. You, "Representation and reinforcement learning for task scheduling in edge computing," *IEEE Transactions on Big Data*, vol. 8, no. 3, pp. 795–808, 2020.
- [111] Q. Qi, L. Zhang, J. Wang, H. Sun, Z. Zhuang, J. Liao, and F. R. Yu, "Scalable parallel task scheduling for autonomous driving using multi-task deep reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 11, pp. 13861–13874, 2020.

- [112] Z. Tong, H. Chen, X. Deng, K. Li, and K. Li, "A scheduling scheme in the cloud computing environment using deep q-learning," *Information Sciences*, vol. 512, pp. 1170–1191, 2020.
- [113] J. Pan and Y. Wei, "A deep reinforcement learning-based scheduling framework for real-time workflows in the cloud environment," *Expert Systems with Applications*, vol. 255, p. 124845, 2024.
- [114] A. M. Kintsakis, F. E. Psomopoulos, and P. A. Mitkas, "Reinforcement learning based scheduling in a workflow management system," *Engineering Applications of Artificial Intelligence*, vol. 81, pp. 94–106, 2019.
- [115] Y. Gu, F. Cheng, L. Yang, J. Xu, X. Chen, and L. Cheng, "Cost-aware cloud workflow scheduling using drl and simulated annealing," *Digital Communications and Networks*, 2024.
- [116] J. Zhang, L. Cheng, C. Liu, Z. Zhao, and Y. Mao, "Cost-aware scheduling systems for real-time workflows in cloud: An approach based on genetic algorithm and deep reinforcement learning," *Expert Systems with Applications*, vol. 234, p. 120972, 2023.
- [117] T. Dong, F. Xue, H. Tang, and C. Xiao, "Deep reinforcement learning for fault-tolerant workflow scheduling in cloud environment," *Applied Intelligence*, vol. 53, no. 9, pp. 9916–9932, 2023.
- [118] H. Li, J. Huang, B. Wang, and Y. Fan, "Weighted double deep q-network based reinforcement learning for bi-objective multi-workflow scheduling in the cloud," *Cluster Computing*, vol. 25, no. 2, pp. 751–768, 2022.
- [119] G. Chen, J. Qi, Y. Sun, X. Hu, Z. Dong, and Y. Sun, "A collaborative scheduling method for cloud computing heterogeneous workflows based on deep reinforcement learning," *Future Generation Computer Systems*, vol. 141, pp. 284–297, 2023.
- [120] S. Zhang, Z. Zhao, C. Liu, and S. Qin, "Data-intensive workflow scheduling strategy based on deep reinforcement learning in multi-clouds," *Journal of Cloud Computing*, vol. 12, no. 1, p. 125, 2023.
- [121] R. Han, S. Wen, C. H. Liu, Y. Yuan, G. Wang, and L. Y. Chen, "EdgeTuner: Fast scheduling algorithm tuning for dynamic edge-cloud workloads and resources," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*, 2022, pp. 880–889.
- [122] T. Zheng, J. Wan, J. Zhang, and C. Jiang, "Deep reinforcement learning-based workload scheduling for edge computing," *Journal of Cloud Computing*, vol. 11, no. 1, p. 3, 2022.
- [123] Y. Qian, L. Shi, J. Li, Z. Wang, H. Guan, F. Shu, and H. V. Poor, "A workflow-aided Internet of things paradigm with intelligent edge computing," *IEEE Network*, vol. 34, no. 6, pp. 92–99, 2020.
- [124] J. Cai, W. Liu, Z. Huang, and F. R. Yu, "Task Decomposition and Hierarchical Scheduling for Collaborative Cloud-Edge-End Computing," *IEEE Transactions on Services Computing*, 2024.
- [125] R. Xie, D. Gu, Q. Tang, T. Huang, and F. R. Yu, "Workflow scheduling in serverless edge computing for the industrial internet of things: A learning approach," *IEEE Transactions on Industrial Informatics*, vol. 19, no. 7, pp. 8242–8252, 2022.
- [126] T. Dong, F. Xue, C. Xiao, and J. Zhang, "Workflow scheduling based on deep reinforcement learning in the cloud environment," *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, no. 12, pp. 10 823–10 835, 2021.
- [127] —, "Deep reinforcement learning for dynamic workflow scheduling in cloud environment," in *2021 IEEE International Conference on Services Computing*. IEEE, 2021, pp. 107–115.
- [128] G. P. Koslovski, K. Pereira, and P. R. Albuquerque, "Dag-based workflows scheduling using actor-critic deep reinforcement learning," *Future Generation Computer Systems*, vol. 150, pp. 354–363, 2024.
- [129] Z. Wang, S. Chen, L. Bai, J. Gao, J. Tao, R. R. Bond, and M. D. Mulvenna, "Reinforcement learning based task scheduling for environmentally sustainable federated cloud computing," *Journal of Cloud Computing*, vol. 12, no. 1, p. 174, 2023.
- [130] J. Xue, T. Wang, and P. Cai, "Towards efficient workflow scheduling over yarn cluster using deep reinforcement learning," in *2023 IEEE Global Communications Conference*, 2023, pp. 473–478.
- [131] H. Peng, C. Wu, Y. Zhan, and Y. Xia, "Lore: a learning-based approach for workflow scheduling in clouds," in *Proceedings of the conference on research in adaptive and convergent systems*, 2022, pp. 47–52.
- [132] L. Lin, L. Pan, and S. Liu, "SpotDAG: An RL-based algorithm for DAG workflow scheduling in heterogeneous cloud environments," *IEEE Transactions on Services Computing*, 2024.
- [133] X. Yang and B. Hu, "An Effective DDPG-generated Task Scheduling Policy to Minimize Latency in Distributed Computing System," in *2023 International Conference on Frontiers of Robotics and Software Engineering*, 2023, pp. 303–310.
- [134] A. Jayanetti, S. Halgamuge, and R. Buyya, "Deep reinforcement learning for energy and time optimized scheduling of precedence-constrained tasks in edge-cloud computing environments," *Future Generation Computer Systems*, vol. 137, pp. 14–30, 2022.
- [135] J. Li, F. Zhou, W. Li, M. Zhao, X. Yan, Y. Xi, and J. Wu, "Componentized Task Scheduling in Cloud-Edge Cooperative Scenarios Based on GNN-enhanced DRL," in *NOMS 2023-2023 IEEE/IFIP Network Operations and Management Symposium*, 2023, pp. 1–4.
- [136] M. Mounesan, M. Lemus, H. Yeddulapalli, P. Calyam, and S. Debroy, "Reinforcement Learning-driven Data-intensive Workflow Scheduling for Volunteer Edge-Cloud," in *2024 IEEE 8th International Conference on Fog and Edge Computing*, 2024, pp. 79–88.
- [137] K. Zhu, Z. Zhang, S. Zeadally, and F. Sun, "Learning to Optimize Workflow Scheduling for an Edge-Cloud Computing Environment," *IEEE Transactions on Cloud Computing*, 2024.
- [138] Z. Wang, M. Goudarzi, M. Gong, and R. Buyya, "Deep reinforcement learning-based scheduling for optimizing system load and response time in edge and fog computing environments," *Future Generation Computer Systems*, vol. 152, pp. 55–69, 2024.
- [139] J. Zhang, T. Wang, and L. Cheng, "Time-Sensitive and Resource-Aware Concurrent Workflow Scheduling for Edge Computing Platforms Based on Deep Reinforcement Learning," *Applied Sciences*, vol. 13, no. 19, p. 10689, 2023.
- [140] K. Zhu, Z. Zhang, F. Sun, and B. Shen, "Workflow makespan minimization for partially connected edge network: A deep reinforcement learning-based approach," *IEEE Open Journal of the Communications Society*, vol. 3, pp. 518–529, 2022.
- [141] A. Asghari, M. K. Sohrabi, and F. Yaghmaee, "A cloud resource management framework for multiple online scientific workflows using cooperative reinforcement learning agents," *Computer Networks*, vol. 179, p. 107340, 2020.
- [142] Y. LI, P. CHEN, K. GUO, and H. XIE, "Multi-Objective Workflow Scheduling With Deep-Q-Network-Based Multi-Agent Reinforcement Learning,"
- [143] A. Jayanetti, S. Halgamuge, and R. Buyya, "Multi-Agent Deep Reinforcement Learning Framework for Renewable Energy-Aware Workflow Scheduling on Distributed Cloud Data Centers," *IEEE Transactions on Parallel and Distributed Systems*, 2024.
- [144] —, "A Deep Reinforcement Learning Approach for Cost Optimized Workflow Scheduling in Cloud Computing Environments," in *Proceedings of the 2024 Asia Pacific Conference on Computing Technologies, Communications and Networking*, 2024, pp. 74–82.
- [145] Y. Duan, J. Li, H. Sun, F. Zhou, J. Chen, T. Wu, W. Li, and Y. Fan, "Telemetry-aided cooperative multi-agent online reinforcement learning for DAG task scheduling in computing power networks," *Simulation Modelling Practice and Theory*, vol. 132, p. 102885, 2024.
- [146] J. Huang, F. Zhou, L. Feng, W. Li, M. Zhao, X. Yan, Y. Xi, and J. Wu, "Digital Twin Assisted DAG Task Scheduling Via Evolutionary Selection MARL in Large-Scale Mobile Edge Network," in *2023 IEEE International Conference on Communications Workshops (ICC Workshops)*, 2023, pp. 158–163.
- [147] Y. Zhao, L. Mo, and J. Liu, "Multi-task scheduling in vehicular edge computing: a multi-agent reinforcement learning approach," *CCF Transactions on Pervasive Computing and Interaction*, pp. 1–17, 2024.
- [148] X. Wang, J. Cao, and R. Buyya, "Adaptive cloud bundle provisioning and multi-workflow scheduling via coalition reinforcement learning," *IEEE Transactions on Computers*, vol. 72, no. 4, pp. 1041–1054, 2022.
- [149] F. Ding, Y. Yuan, L. Lv, R. Zhang, and W. Zhou, "Transformer-Enhanced DQN Approach for Energy and Cost-Efficient Large-Scale Dynamic Workflow Scheduling in Heterogeneous Environment," *IEEE Internet of Things Journal*, 2024.
- [150] Z. Liu, L. Huang, Z. Gao, M. Luo, S. Hosseinalipour, and H. Dai, "GA-DRL: Graph Neural Network-Augmented Deep Reinforcement Learning for DAG Task Scheduling over Dynamic Vehicular Clouds," *IEEE Transactions on Network and Service Management*, 2024.
- [151] T. Long, Y. Xia, Y. Ma, Q. Peng, and J. Zhao, "A fault-tolerant workflow scheduling method on deep reinforcement learning-based in edge environment," in *2022 IEEE International Conference on Networking, Sensing and Control*, 2022, pp. 1–6.
- [152] A. Mahapatra, R. Pradhan, S. K. Majhi, and K. Mishra, "Quantum ML-Based Cooperative Task Orchestration in Dew-Assisted IoT Framework," *Arabian Journal for Science and Engineering*, pp. 1–28, 2024.
- [153] Z. Wang, M. Goudarzi, and R. Buyya, "TF-DDRL: A Transformer-enhanced Distributed DRL Technique for Scheduling IoT Applications in Edge and Cloud Computing Environments," *arXiv preprint arXiv:2410.14348*, 2024.

- [154] E. Deelman, K. Vahi, G. Juve, M. Rynge, S. Callaghan, P. J. Maechling, R. Mayani, W. Chen, R. F. Da Silva, M. Livny *et al.*, "Pegasus, a workflow management system for science automation," *Future Generation Computer Systems*, vol. 46, pp. 17–35, 2015.
- [155] Q. Zong, X. Zheng, Y. Wei, and H. Sun, "A deep reinforcement learning based resource autonomic provisioning approach for cloud services," in *Collaborative Computing: Networking, Applications and Worksharing*. Cham: Springer International Publishing, 2021, pp. 132–153.
- [156] S. Tuli, G. Casale, and N. R. Jennings, "CILP: Co-Simulation-Based Imitation Learner for Dynamic Resource Provisioning in Cloud Computing Environments," *IEEE Transactions on Network and Service Management*, vol. 20, no. 4, pp. 4448–4460, 2023.
- [157] M. Zhu, Q. Chen, J. Gu, and P. Gu, "Deep reinforcement learning for provisioning virtualized network function in inter-datacenter elastic optical networks," *IEEE Transactions on Network and Service Management*, vol. 19, no. 3, pp. 3341–3351, 2022.
- [158] H. Sami, H. Otrok, J. Bentahar, and A. Mourad, "AI-based resource provisioning of IoE services in 6G: A deep reinforcement learning approach," *IEEE Transactions on Network and Service Management*, vol. 18, no. 3, pp. 3527–3540, 2021.
- [159] R. Zhu, G. Li, P. Wang, M. Xu, and S. Yu, "Drl-based deadline-driven advance reservation allocation in eons for cloud-edge computing," *IEEE Internet of Things Journal*, vol. 9, no. 21, pp. 21 444–21 457, 2022.
- [160] M. Faraji-Mehmandar, S. Jabbehdari, and H. H. S. Javadi, "A self-learning approach for proactive resource and service provisioning in fog environment," *The Journal of Supercomputing*, vol. 78, no. 15, pp. 16 997–17 026, 2022.
- [161] J. Santos, T. Wauters, B. Volckaert, and F. D. Turck, "Resource Provisioning in Fog Computing through Deep Reinforcement Learning," 2021.
- [162] Z. Chen, J. Hu, G. Min, C. Luo, and T. El-Ghazawi, "Adaptive and efficient resource allocation in cloud datacenters using actor-critic deep reinforcement learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 8, pp. 1911–1923, 2021.
- [163] W. Funika, P. Koperek, and J. Kitowski, "Automated cloud resources provisioning with the use of the proximal policy optimization," *The Journal of Supercomputing*, vol. 79, no. 6, pp. 6674–6704, 2023.
- [164] X. Zhang, B. Li, J. Peng, X. Pan, and Z. Zhu, "You calculate and I provision: A DRL-assisted service framework to realize distributed and tenant-driven virtual network slicing," *Journal of Lightwave Technology*, vol. 39, no. 1, pp. 4–16, 2021.
- [165] H. Baghban, A. Rezapour, C.-H. Hsu, S. Nuannimnoi, and C.-Y. Huang, "Edge-AI: IoT Request Service Provisioning in Federated Edge Computing Using Actor-Critic Reinforcement Learning," *IEEE Transactions on Engineering Management*, vol. 71, pp. 12 519–12 528, 2024.
- [166] S. Guo, Y. Dai, S. Xu, X. Qiu, and F. Qi, "Trusted cloud-edge network resource management: DRL-driven service function chain orchestration for IoT," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6010–6022, 2019.
- [167] A. S. A. Geetha, and R. K., "Intelligent Resource Provisioning and Optimization in Fog Computing using Deep Reinforcement Learning," *International Journal of Electronics and Communication Engineering*, vol. 10, no. 8, pp. 85–97, 2023.
- [168] M. Chen, T. Wang, S. Zhang, and A. Liu, "Deep reinforcement learning for computation offloading in mobile edge computing environment," *Computer Communications*, vol. 175, pp. 1–12, 2021.
- [169] J. Zhang, S. Chen, X. Wang, and Y. Zhu, "Dynamic reservation of edge servers via deep reinforcement learning for connected vehicles," *IEEE Transactions on Mobile Computing*, vol. 22, no. 5, pp. 2661–2674, 2021.
- [170] A. Jyoti and M. Shrimali, "Dynamic provisioning of resources based on load balancing and service broker policy in cloud computing," *Cluster Computing*, vol. 23, no. 1, pp. 377–395, 2020.
- [171] A. Asghari and M. K. Sohrabi, "Combined use of coral reefs optimization and multi-agent deep Q-network for energy-aware resource provisioning in cloud data centers using DVFS technique," *Cluster Computing*, vol. 25, no. 1, pp. 119–140, 2022.
- [172] Y. Zhang, B. Di, Z. Zheng, J. Lin, and L. Song, "Distributed Multi-Cloud Multi-Access Edge Computing by Multi-Agent Reinforcement Learning," *IEEE Transactions on Wireless Communications*, vol. 20, no. 4, pp. 2565–2578, 2021.
- [173] H. Li, Y. Liu, X. Zhou, X. Vasilakos, R. Nejabati, S. Yan, and D. Simeonidou, "Adaptive resource management for edge network slicing using incremental multi-agent deep reinforcement learning," *arXiv preprint arXiv:2310.17523*, 2023.
- [174] M. Asim Ejaz, G. Wu, and T. Iqbal, "Dynamic and efficient resource allocation for 5G end-to-end network slicing: A multi-agent deep reinforcement learning approach," *International Journal of Communication Systems*, p. e5916.
- [175] A. Asghari, M. K. Sohrabi, and F. Yaghmaee, "Task scheduling, resource provisioning, and load balancing on scientific workflows using parallel SARSA reinforcement learning agents and genetic algorithm," *The Journal of Supercomputing*, vol. 77, no. 3, pp. 2800–2828, 2021.
- [176] M. Xu, D. Niyato, J. Kang, Z. Xiong, Y. Cao, Y. Gao, C. Ren, and H. Yu, "Generative AI-enabled Quantum Computing Networks and Intelligent Resource Allocation," 2024.
- [177] A. Shahidinejad and M. Ghobaei-Arani, "Joint computation offloading and resource provisioning for edge-cloud computing environment: A machine learning-based approach," *Software: Practice and Experience*, vol. 50, no. 12, pp. 2212–2230, 2020.
- [178] B. He, J. Wang, Q. Qi, H. Sun, and J. Liao, "Towards Intelligent Provisioning of Virtualized Network Functions in Cloud of Things: A Deep Reinforcement Learning Based Approach," *IEEE Transactions on Cloud Computing*, vol. 10, no. 2, pp. 1262–1274, 2022.
- [179] M. Dieye, W. Jaafar, H. Elbiaze, and R. H. Glitho, "DRL-Based Green Resource Provisioning for 5G and Beyond Networks," *IEEE Transactions on Green Communications and Networking*, vol. 7, no. 4, pp. 2163–2180, 2023.
- [180] J. Chen, S. Chen, Q. Wang, B. Cao, G. Feng, and J. Hu, "iRAF: A deep reinforcement learning approach for collaborative mobile edge computing IoT networks," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 7011–7024, 2019.
- [181] S. Kardani-Moghaddam, R. Buyya, and K. Ramamohanarao, "ADRL: A Hybrid Anomaly-Aware Deep Reinforcement Learning-Based Resource Scaling in Clouds," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 3, pp. 514–526, 2021.
- [182] X. Chen, L. Yang, Z. Chen, G. Min, X. Zheng, and C. Rong, "Resource Allocation With Workload-Time Windows for Cloud-Based Software Services: A Deep Reinforcement Learning Approach," *IEEE Transactions on Cloud Computing*, vol. 11, no. 2, pp. 1871–1885, 2023.
- [183] C. Bitsakos, I. Konstantinou, and N. Koziris, "DERP: A Deep Reinforcement Learning Cloud System for Elastic Resource Provisioning," in *2018 IEEE International Conference on Cloud Computing Technology and Science*, 2018, pp. 21–29.
- [184] T. Cui, R. Yang, C. Fang, and S. Yu, "Deep Reinforcement Learning-Based Resource Allocation for Content Distribution in IoT-Edge-Cloud Computing Environments," *Symmetry*, vol. 15, no. 1, p. 217, 2023.
- [185] X. Li, "A Computing Offloading Resource Allocation Scheme Using Deep Reinforcement Learning in Mobile Edge Computing Systems," *Journal of Grid Computing*, vol. 19, no. 3, p. 35, 2021.
- [186] C. Fang, H. Xu, Y. Yang, Z. Hu, S. Tu, K. Ota, Z. Yang, M. Dong, Z. Han, F. R. Yu, and Y. Liu, "Deep-Reinforcement-Learning-Based Resource Allocation for Content Distribution in Fog Radio Access Networks," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 16 874–16 883, 2022.
- [187] H. Zhang, J. Wang, H. Zhang, and C. Bu, "Security computing resource allocation based on deep reinforcement learning in serverless multi-cloud edge computing," *Future Generation Computer Systems*, vol. 151, pp. 152–161, 2024.
- [188] I. Ullah, H.-K. Lim, Y.-J. Seok, and Y.-H. Han, "Optimizing task offloading and resource allocation in edge-cloud networks: A DRL approach," *Journal of Cloud Computing*, vol. 12, no. 1, p. 112, 2023.
- [189] T. Liu, S. Ni, X. Li, Y. Zhu, L. Kong, and Y. Yang, "Deep Reinforcement Learning Based Approach for Online Service Placement and Computation Resource Allocation in Edge Computing," *IEEE Transactions on Mobile Computing*, vol. 22, no. 7, pp. 3870–3881, 2023.
- [190] H. Mao, M. Alizadeh, I. Menache, and S. Kandula, "Resource Management with Deep Reinforcement Learning," in *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*, 2016, pp. 50–56.
- [191] M. Arvindhan and D. R. Kumar, "Adaptive Resource Allocation in Cloud Data Centers using Actor-Critical Deep Reinforcement Learning for Optimized Load Balancing," *International Journal on Recent and Innovation Trends in Computing and Communication*, vol. 11, no. 5s, pp. 310–318, 2023.
- [192] Z. Chen, J. Hu, and G. Min, "Learning-Based Resource Allocation in Cloud Data Center using Advantage Actor-Critic," in *2019 IEEE International Conference on Communications (ICC)*, 2019, pp. 1–6.
- [193] Y. He, Y. Wang, C. Qiu, Q. Lin, J. Li, and Z. Ming, "Blockchain-Based Edge Computing Resource Allocation in IoT: A Deep Reinforcement Learning Approach," *IEEE Internet of Things Journal*, vol. 8, no. 4, pp. 2226–2237, 2021.

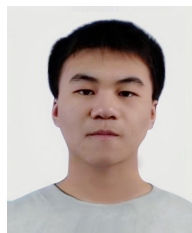
- [194] X. Zhu, Y. Luo, A. Liu, N. N. Xiong, M. Dong, and S. Zhang, "A Deep Reinforcement Learning-Based Resource Management Game in Vehicular Edge Computing," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 3, pp. 2422–2433, 2022.
- [195] J. Huang, J. Wan, B. Lv, Q. Ye, and Y. Chen, "Joint Computation Offloading and Resource Allocation for Edge-Cloud Collaboration in Internet of Vehicles via Deep Reinforcement Learning," *IEEE Systems Journal*, vol. 17, no. 2, pp. 2500–2511, 2023.
- [196] A. Belgacem, S. Mahmoudi, and M. Kihl, "Intelligent multi-agent reinforcement learning model for resources allocation in cloud computing," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 6, pp. 2391–2404, 2022.
- [197] S. Nagarajan, P. S. Rani, M. S. Vinmathi, V. Subba Reddy, A. L. M. Saleth, and D. Abdus Subhahan, "Multi agent deep reinforcement learning for resource allocation in container-based clouds environments," *Expert Systems*, p. exsy.13362, 2023.
- [198] X. Liu, J. Yu, Z. Feng, and Y. Gao, "Multi-agent reinforcement learning for resource allocation in IoT networks with edge computing," *China Communications*, vol. 17, no. 9, pp. 220–236, 2020.
- [199] J. Rosenberger, M. Urlaub, F. Rauterberg, T. Lutz, A. Selig, M. Bühren, and D. Schramm, "Deep Reinforcement Learning Multi-Agent System for Resource Allocation in Industrial Internet of Things," *Sensors*, vol. 22, no. 11, p. 4099, 2022.
- [200] W. Guo, W. Tian, Y. Ye, L. Xu, and K. Wu, "Cloud Resource Scheduling With Deep Reinforcement Learning and Imitation Learning," *IEEE Internet of Things Journal*, vol. 8, no. 5, pp. 3576–3586, 2021.
- [201] Y. Zhang, J. Yao, and H. Guan, "Intelligent Cloud Resource Management with Deep Reinforcement Learning," *IEEE Cloud Computing*, vol. 4, no. 6, pp. 60–69, 2017.
- [202] W. Wei, H. Gu, K. Wang, J. Li, X. Zhang, and N. Wang, "Multi-Dimensional Resource Allocation in Distributed Data Centers Using Deep Reinforcement Learning," *IEEE Transactions on Network and Service Management*, vol. 20, no. 2, pp. 1817–1829, 2023.
- [203] F. Jiang, K. Wang, L. Dong, C. Pan, and K. Yang, "Stacked Autoencoder-Based Deep Reinforcement Learning for Online Resource Scheduling in Large-Scale MEC Networks," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9278–9290, 2020.
- [204] F. Xue, Q. Hai, T. Dong, Z. Cui, and Y. Gong, "A deep reinforcement learning based hybrid algorithm for efficient resource scheduling in edge computing environment," *Information Sciences*, vol. 608, pp. 362–374, 2022.
- [205] J. Adu Ansere, D. T. Tran, O. A. Dobre, H. Shin, G. K. Karagiannis, and T. Q. Duong, "Energy-Efficient Optimization for Mobile Edge Computing With Quantum Machine Learning," *IEEE Wireless Communications Letters*, vol. 13, no. 3, pp. 661–665, 2024.
- [206] Silvirianti, B. Narottama, and S. Y. Shin, "Layerwise Quantum Deep Reinforcement Learning for Joint Optimization of UAV Trajectory and Resource Allocation," *IEEE Internet of Things Journal*, vol. 11, no. 1, pp. 430–443, 2024.
- [207] W. Funika, P. Koperek, and J. Kitowski, "Management of heterogeneous cloud resources with use of the ppo," in *European Conference on Parallel Processing*. Springer, 2020, pp. 148–159.
- [208] Z. Sun, H. Huang, Z. Li, C. Gu, R. Xie, and B. Qian, "Efficient, economical and energy-saving multi-workflow scheduling in hybrid cloud," *Expert Systems with Applications*, vol. 228, p. 120401, 2023.
- [209] Z. Sun, B. Zhang, C. Gu, R. Xie, B. Qian, and H. Huang, "Et2fa: A hybrid heuristic algorithm for deadline-constrained workflow scheduling in cloud," *IEEE Transactions on Services Computing*, vol. 16, no. 3, pp. 1807–1821, 2022.
- [210] Y. Wen, J. Liu, W. Dou, X. Xu, B. Cao, and J. Chen, "Scheduling workflows with privacy protection constraints for big data applications on cloud," *Future Generation Computer Systems*, vol. 108, pp. 1084–1091, 2020.
- [211] M. N. Alenezi, H. Alabdulrazzaq, and N. Q. Mohammad, "Symmetric encryption algorithms: Review and evaluation study," *International Journal of Communication Networks and Information Security*, vol. 12, no. 2, pp. 256–272, 2020.
- [212] S. R. Prasanna and B. Premananda, "Performance analysis of md5 and sha-256 algorithms to maintain data integrity," in *2021 International Conference on Recent Trends on Electronics, Information, Communication & Technology*, 2021, pp. 246–250.
- [213] C. Zhou, J. Gao, M. Li, X. S. Shen, and W. Zhuang, "Digital twin-empowered network planning for multi-tier computing," *Journal of Communications and Information Networks*, vol. 7, no. 3, pp. 221–238, 2022.
- [214] S. O. Oladejo, S. O. Ekwe, and L. A. Akinyemi, "Multi-tier multi-domain network slicing: A resource allocation perspective," in *2021 IEEE AFRICON*, 2021, pp. 1–6.
- [215] R. Ou, G. Sun, D. Ayepah-Mensah, G. O. Boateng, and G. Liu, "Two-Tier Resource Allocation for Multitenant Network Slicing: A Federated Deep Reinforcement Learning Approach," *IEEE Internet of Things Journal*, vol. 10, no. 22, pp. 20174–20187, 2023.
- [216] Y. Guan, Y. Liu, Y. Li, and X. Xu, "HierRL: Hierarchical reinforcement learning for task scheduling in distributed systems," in *2022 International Joint Conference on Neural Networks*, 2022, pp. 1–8.
- [217] J. Zhang, B. Guo, X. Ding, D. Hu, J. Tang, K. Du, C. Tang, and Y. Jiang, "An adaptive multi-objective multi-task scheduling method by hierarchical deep reinforcement learning," *Applied Soft Computing*, vol. 154, p. 111342, 2024.
- [218] C. Ma, A. Li, Y. Du, H. Dong, and Y. Yang, "Efficient and scalable reinforcement learning for large-scale network control," *Nature Machine Intelligence*, pp. 1–15, 2024.
- [219] S. Lai, Z. Xu, W. Zhang, H. Liu, and H. Xiong, "Large language models as traffic signal control agents: Capacity and opportunity," *arXiv preprint arXiv:2312.16044*, 2023.



Yan Gu received the B.E. degree from Nanjing Institute of Technology, Nanjing, China, in 2020, and M.S. degree from the School of Computer at Jiangsu University of Science and Technology, China. She is currently a PhD student in the School of Control and Computer Engineering at North China Electric Power University in Beijing. Her research interests include cloud computing, deep learning, parallel and distributed computing.



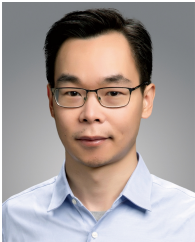
Zhaoze Liu received the B.E. degree from Beijing University of Posts and Telecommunications, Beijing, China, in 2022. He is currently pursuing the M.S. degree in the School of Control and Computer Engineering at North China Electric Power University in Beijing. His research interests include cloud computing, deep reinforcement learning, and serverless computing.



Shuhong Dai is a master student in Software Engineering at North China Electric Power University in Beijing. He received his B.E. degree in Electrical Engineering and Automation from Ningbo University of Technology, China, in 2023. His research interests include deep reinforcement learning, cloud computing and intelligent transportation systems.



Cong Liu received the PhD degree in the Department of Mathematics and Computer Science, Eindhoven University of Technology in 2019. He is a full Professor in Shandong University of Technology. His research interests are in the areas of business process management, process mining, and workflow management.



Ying Wang is a Professor at Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS). He received Ph.D degree from ICT in 2014. His research interests includes computer architecture and VLSI design, specifically memory system, on-chip interconnects, resilient and energy-efficient architecture, machine learning accelerators, and parallel data systems.



Shen Wang is an Assistant Professor with the School of Computer Science, University College Dublin, Ireland. He received the Ph.D. degree from Dublin City University, Ireland. Dr. Wang has been involved with several EU projects as a co-PI, WP and Task leader in big trajectory data streaming for air traffic control and trustworthy AI for intelligent cybersecurity systems. His research interests include connected autonomous vehicles, deep reinforcement learning, and security and privacy for mobile networks.



Georgios Theodoropoulos is currently a Chair Professor at the Department of Computer Science and Engineering at SUSTech in Shenzhen, China. He was previously the inaugural Executive Director of the Institute of Advanced Research Computing and a Chair Professor at the School of Engineering and Computing Sciences at the University of Durham, UK. He has been a Senior Research Scientist with IBM Research and senior faculty at the University of Birmingham, UK, where he was also founding Director of one of UK's e-Science Centres of Excellence. He has held an Adjunct Chair at the Trinity College Dublin and visiting appointments at the Nanyang Technological University and National University in Singapore. He is Ordinary Member of the European Academy of Sciences and Arts, a Fellow of the World Academy of Art and Science, an Accredited Board Director of the Singapore Institute of Directors, a Chartered Engineer, and holds a PhD from the University of Manchester, UK.



Long Cheng is a Professor at North China Electric Power University in Beijing. He received the Ph.D from National University of Ireland Maynooth in 2014. He was an Assistant Professor at Dublin City University, and a Marie Curie Fellow at University College Dublin. He has published more than 110 papers in refereed journals and conferences. His research focuses on distributed computing and deep reinforcement learning. Prof Cheng is an Associate Editor of IEEE Transactions on Consumer Electronics, and a Chair of Journal of Cloud Computing.

More info: <https://longcheng.eu/>